

Міністерство освіти і науки України
Сумський державний університет
Навчально-науковий інститут бізнес-технологій «УАБС»
Кафедра економічної кібернетики

КВАЛІФІКАЦІЙНА МАГІСТЕРСЬКА РОБОТА

на тему «ІНТЕЛЕКТУАЛЬНИЙ АНАЛІЗ КІБЕРЗАГРОЗ В БАНКАХ»

Виконав студент 2 курсу, групи ЕК.м-61а

(номер курсу)

(шифр групи)

Спеціальності 051 «Економіка («Економічна
кібернетика»))»

Сковронська Анастасія Ігорівна

(прізвище, ініціали студента)

Керівник к.е.н., доцент, Яровенко Г. М.

(посада, науковий ступінь, прізвище, ініціали)

Суми – 2018 рік

РЕФЕРАТ

кваліфікаційної магістерської роботи на тему «ІНТЕЛЕКТУАЛЬНИЙ АНАЛІЗ КІБЕРЗАГРОЗ В БАНКАХ»

студента Сковронської Анастасії Ігорівни
(прізвище, ім'я, по батькові)

Актуальність кваліфікаційної магістерської роботи визначається високим рівнем уразливості банківської сфери до кіберзагроз, значними витратами на кібербезпеку у банку, швидкими темпами розвитку нових кібернетичних загроз.

Мета кваліфікаційної магістерської роботи полягає у побудові математичної моделі для виявлення ознак кіберзагроз у банку та її практичній реалізації із використанням методів інтелектуального аналізу за допомогою аналітичного пакету SAS Enterprise Miner.

Об'єктом дослідження є взаємовідносини учасників банківської діяльності, в результаті яких створюються умови для виникнення кіберзагроз, що становлять небезпеку банківській сфері

Предметом дослідження виступають методи та моделі інтелектуального аналізу для виявлення ознак кіберзагроз у банку.

Методи дослідження: аналіз і синтез, дедукція, абстрагування, конкретизація, аргументація, порівняння, класифікація та узагальнення.

Інформаційна база кваліфікаційної магістерської роботи: нормативно-довідкова та рекомендаційна документація компанії SAS, емпіричні дані банку по транзакціям користувачів мобільного та інтернет-банкінгу.

В результаті роботи сформульовані наступні висновки: серед побудованих моделей за показниками якості та адекватності найкращою було обрано нейронну мережу та виявлено, що на тестовій вибірці, 21,9% транзакцій ймовірно є кіберзагрозами.

Отриману модель можливо застосовувати в банках для виявлення та попередження можливого виникнення ознак кіберзагроз в транзакціях користувачів мобільного та інтернет-банкінгу.

Ключові слова: кіберзагроза, банківська транзакція, інтелектуальний аналіз, моделювання, кластерний аналіз, регресійна модель, дерево рішень, нейронна мережа, SAS Enterprise Miner.

Результати проведеної роботи було апробовано та опубліковано в тезах II-ї Міжнародної науково-практичної конференції «Теорія і практика розвитку наукових знань», що відбувалась в м. Київ 28-29 грудня 2017 року.

Основний зміст кваліфікаційної магістерської роботи викладено на 101 сторінці, зокрема список використаних джерел із 40 найменувань, розміщений на 4 сторінках. Робота містить 16 таблиць, 48 рисунків.

Рік виконання кваліфікаційної роботи 2017 – 2018.

Рік захисту роботи 2018.

Міністерство освіти і науки України
Сумський державний університет
Навчально-науковий інститут бізнес-технологій «УАБС»
Кафедра економічної кібернетики

ЗАТВЕРДЖУЮ
Завідувач кафедри
д.е.н., доцент
_____ О.В. Кузьменко
«__» _____ 2017 р.

ЗАВДАННЯ НА КВАЛІФІКАЦІЙНУ МАГІСТЕРСЬКУ РОБОТУ
(спеціальність 051 «Економіка («Економічна кібернетика»))»
студенту 2 курсу, групи ЕК.м-61а

Сковронської Анастасії Ігорівни
(прізвище, ім'я, по батькові студента)

1. Тема роботи Інтелектуальний аналіз кіберзагроз в банках
затверджена наказом по університету від «08» грудня 2017 року № 2748-III
2. Термін подання студентом закінченої роботи « » _____ 2018 року
3. Мета кваліфікаційної роботи полягає у побудові математичної моделі для виявлення ознак кіберзагроз у банку та її практичній реалізації із використанням методів інтелектуального аналізу за допомогою аналітичного пакету SAS Enterprise Miner.
4. Об'єкт дослідження: взаємовідносини учасників банківської діяльності, в результаті яких створюються умови для виникнення кіберзагроз, що становлять небезпеку банківській сфері.
5. Предмет дослідження: методи та моделі інтелектуального аналізу для виявлення ознак кіберзагроз у банку.
6. Кваліфікаційна робота виконується на матеріалах: нормативно-довідкова та рекомендаційна документація компанії SAS, емпіричні дані банку по транзакціям користувачів мобільного та інтернет-банкінгу.

7. Орієнтовний план кваліфікаційної роботи, терміни подання розділів керівникові та зміст завдань для виконання поставленої мети

Розділ 1 Теоретико-методологічні основи інтелектуального аналізу кіберзагроз в банках – 13 листопада 2017 р.

(назва – термін подання)

У розділі 1 розкрити сутність кіберзагроз як об'єкту моделювання, проаналізувати існуючі підходи до виявлення кіберзагроз у банках, дослідити методи інтелектуального аналізу, побудувати концептуальну модель виявлення ознак кіберзагроз

(зміст конкретних завдань до розділу, які повинен виконати студент)

Розділ 2 Математичні моделі виявлення ознак кіберзагроз у банках – 12 грудня 2017 р.

(назва – термін подання)

У розділі 2 провести первинний аналіз вхідних та вихідних даних, розглянути кластерний аналіз як інструмент дослідження первинних даних, обґрунтувати вибір методів інтелектуального аналізу для виявлення ознак кіберзагроз

(зміст конкретних завдань до розділу, які має виконати студент)

Розділ 3 Практична реалізація інтелектуального аналізу ознак кіберзагроз із використанням аналітичного пакету SAS Enterprise Miner – 4 січня 2018 р.

(назва – термін подання)

У розділі 3 провести опис програмної реалізації модельних розрахунків на ЕОМ, проаналізувати якість та адекватність побудованих моделей та порівняти їх, провести оцінку результатів та ефекту від застосування моделі

(зміст конкретних завдань до розділу, які повинен виконати студент)

8. Консультації з роботи:

Розділ	Прізвище, ініціали та посада консультанта	Підпис, дата	
		завдання видав	завдання прийняв
1		Яровенко Г. М.	Сковронська А. І.
2		Яровенко Г. М.	Сковронська А. І.
3		Яровенко Г. М.	Сковронська А. І.

9. Дата видачі завдання: «___» _____ 20__ року

Керівник кваліфікаційної роботи

(підпис)

Г. М. Яровенко
(ініціали, прізвище)

Завдання до виконання одержав

(підпис)

А. І. Сковронська
(ініціали, прізвище)

ЗМІСТ

ВСТУП	7
РОЗДІЛ 1 ТЕОРЕТИКО-МЕТОДОЛОГІЧНІ ОСНОВИ ІНТЕЛЕКТУАЛЬНОГО АНАЛІЗУ КІБЕРЗАГРОЗ В БАНКАХ	10
1.1 Аналіз кіберзагроз як об'єкту моделювання.....	10
1.2 Аналіз існуючих підходів до виявлення кіберзагроз у банках	20
1.3 Аналіз методів інтелектуального аналізу	24
1.4 Побудова концептуальної моделі виявлення ознак кіберзагроз.....	29
РОЗДІЛ 2 МАТЕМАТИЧНІ МОДЕЛІ ВИЯВЛЕННЯ ОЗНАК КІБЕРЗАГРОЗ У БАНКАХ.....	34
2.1 Первинний аналіз вхідних та вихідних даних	34
2.2 Кластерний аналіз як інструмент дослідження первинних даних ...	38
2.3 Обґрунтування вибору методів інтелектуального аналізу для виявлення ознак кіберзагроз у банку	42
РОЗДІЛ 3 ПРАКТИЧНА РЕАЛІЗАЦІЯ ІНТЕЛЕКТУАЛЬНОГО АНАЛІЗУ ОЗНАК КІБЕРЗАГРОЗ ІЗ ВИКОРИСТАННЯМ АНАЛІТИЧНОГО ПАКЕТУ SAS ENTERPRISE MINER.....	57
3.1 Опис програмної реалізації модельних розрахунків на ЕОМ.....	57
3.2 Аналіз якості та адекватності побудованих моделей.....	79
3.3 Оцінка результатів та ефекту.....	88
ВИСНОВКИ.....	94
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ	97
ДОДАТКИ.....	101

ВСТУП

В сучасній ері інформаційного суспільства комп'ютери та телекомунікаційні системи охоплюють всі сфери життя людей і країн. Але людство, використовуючи телекомунікації та глобальні комп'ютерні мережі, не уявляє, які можливості для зловживання створюють дані технології.

Кіберзагрози в сучасному суспільстві набирають значного масштабу. Відтепер успішна атака хакерів може знеструмити область або країну, призвести до пограбування банку чи знищити успішну організацію.

За результатами опитування, представленого компанією IBM, кожен випадок хакерської атаки чи витоку даних обійшовся банкам і фінансовим установам в середньому в 3,79 мільйонів доларів. Чверть всіх подібних інцидентів були викликані людським фактором. В переважній більшості випадків співробітники банків перейшли по зловмисним посиланням, відкривали вкладення в фішингових повідомленнях або ставали жертвами складних атак з використанням методів соціального інжинірингу. 29% проблем припали на долю збитків у роботі обладнання та програмного забезпечення самих фінансових організацій. Однак більша частина інцидентів (47%) пов'язана з DDoS- та хакерськими атаками безпосередньо на системи самих банків.

Актуальність даної проблеми визначається високим рівнем уразливості банківської сфери до кіберзагроз, значними витратами на кібербезпеку у банку, швидкими темпами розвитку нових кібернетичних загроз.

Банківський сектор займає верхні рядки списку цільових галузей, найбільш схильних до атак кіберзлочинців. Саме тому банки повинні завжди бути пильними і оперативно реагувати на інформацію про потенційні загрози, типи атак і методи їх реалізації.

Об'єктом дослідження є взаємовідносини учасників банківської діяльності, в результаті яких створюються умови для виникнення кіберзагроз, що становлять небезпеку банківській сфері.

Предметом дослідження виступають методи та моделі інтелектуального аналізу для виявлення ознак кіберзагроз у банку.

Обрані об'єкт та предмет обумовили мету дослідження, яка полягає у побудові математичної моделі для виявлення ознак кіберзагроз у банку та її практичній реалізації із використанням методів інтелектуального аналізу за допомогою аналітичного пакету SAS Enterprise Miner.

Для вирішення поставленої мети були визначені наступні завдання:

- а) розкрити сутність об'єкту моделювання – взаємовідносин учасників банківської діяльності, в результаті яких створюються умови для виникнення кіберзагроз, що становлять небезпеку банківській сфері;
- б) проаналізувати існуючі підходи до виявлення кіберзагроз у банках;
- в) дослідити методи інтелектуального аналізу;
- г) побудувати концептуальну модель виявлення ознак кіберзагроз;
- д) провести первинний аналіз вхідних та вихідних даних;
- е) застосувати кластерний аналіз як інструмент дослідження первинних даних;
- ж) обґрунтувати вибір методів інтелектуального аналізу для виявлення ознак кіберзагроз у банку;
- з) описати програмну реалізацію модельних розрахунків із використанням аналітичного пакету SAS Enterprise Miner;
- и) проаналізувати якість та адекватність побудованих моделей;
- к) провести оцінку результатів та ефекту.

При дослідженні теми в роботі було використано такі загальнонаукові методи: аналіз і синтез, дедукція, абстрагування, конкретизація, аргументація, порівняння, класифікація та метод узагальнення, за допомогою якого було зроблено загальні висновки.

Інформаційно-фактологічну базу склали: нормативно-довідкова та рекомендаційна документація компанії SAS; емпіричні дані банку по транзакціям користувачів мобільного та інтернет-банкінгу, наукові публікації вчених-фахівців в галузі економіко-математичного моделювання; нормативно-правова документація Національного банку України та інших банків країни.

Результати проведеної роботи було апробовано та опубліковано в тезах II-ї Міжнародної науково-практичної конференції «Теорія і практика розвитку наукових знань», що відбувалась в м. Київ 28-29 грудня 2017 року.

РОЗДІЛ 1 ТЕОРЕТИКО-МЕТОДОЛОГІЧНІ ОСНОВИ ІНТЕЛЕКТУАЛЬНОГО АНАЛІЗУ КІБЕРЗАГРОЗ В БАНКАХ

1.1 Аналіз кіберзагроз як об'єкту моделювання

Щоденна діяльність банківських систем тісно пов'язана з використанням сучасних комп'ютерних технологій і перебуває в повній залежності від надійної та безперебійної роботи електронно-обчислювальних систем. Світовий досвід свідчить про безумовну уразливість будь-якої компанії з огляду на те, що кіберзлочини не мають державних кордонів, у зв'язку з чим хакери мають можливість в рівній мірі загрожувати інформаційним системам в будь-якій точці світу [24].

Кібернетична загроза (кіберзагроза) – наявні й потенційно можливі явища та чинники, що створюють небезпеку інтересам людини, суспільства й держави через порушення доступності, повноти, цілісності, достовірності, автентичності режиму доступу до інформації, яка циркулює в критичних об'єктах національної інформаційної інфраструктури [36].

Основоположні причини виникнення кіберзагроз полягають в:

- відсутності необхідного законодавства і єдиних стандартів безпеки;
- недостатності фінансування з боку самих банків;
- відсутності корпоративної культури в сфері кібербезпеки всередині банку [24].

Розглянемо найпоширеніші кіберзагрози в банках:

а) атаки мережевого та прикладного рівнів:

- 1) розрив або призупинення серверів та мережевих ресурсів, підключених до Інтернету;
- 2) легка атака для будь-кого для запуску, дуже важко для банків вирішити самотійно;

- 3) пакети атак DDoS легко доступні будь-кому на чорному ринку;
- 4) атаки DDoS можуть запускатися кіберзлочинцями, щоб відвернути банківський персонал від помітних шахрайських операцій, таких як несанкціоновані перекази коштів;

б) соціальна інженерія:

- 1) банківські клієнти часто натрапляють на фішингові атаки;
- 2) банківські клієнти отримують підроблені електронні листи, які використовуються для отримання доступу до їх рахунків або отримання особистої інформації;
- 3) підроблені електронні листи ретельно створюються, щоб відобразити справжні листи, які зазвичай надсилаються банками.
- 4) важко виявити, оскільки джерело електронної пошти часто виявляється законним.

в) розвинені стійкі загрози:

- 1) «Backdoor» для систем встановлюється за допомогою вразливостей («Backdoor» - вразливість в програмі, що дозволяє хакерам зламати систему або здійснити будь-яку недружелюбну дію);
- 2) за допомогою належного шкідливого коду нападники залишаються непоміченими, щоб як можна довше продовжувати наносити збитки;

г) організована кіберзлочинність:

- 1) ризик розкрадання інтелектуальної власності, конфіскація банківських рахунків та втрата споживачів внаслідок бізнес-збоїв;
- 2) в кінцевому рахунку, легше запобігти, ніж усунути, кіберзлочинці спеціалізуються на продажі особистої інформації на чорному ринку, використовуючи викуп та шантаж;

д) порушення основних даних:

- 1) високоорганізовані хакери, які використовують надійну інфраструктуру для цільових банківських установ, викрадають дані

клієнтів та продають їх;

2) за допомогою різних методів розкривається конфіденційна інформація про банківські установи та їх клієнтів;

3) бізнес порушується, дані про клієнтів та компанії погіршуються, а витрати на відновлення є величезними [12].

DoS (від англ. Denial of Service – відмова в обслуговуванні) – хакерська атака на обчислювальну систему з метою довести її до відмови, тобто створення таких умов, при яких сумлінні користувачі системи не можуть отримати доступ до надаваних системних ресурсів (серверів), або цей доступ ускладнений. Відмова «ворожої» системи може бути і кроком до оволодіння системою. Але частіше – це міра економічного тиску: втрата звичайної служби, що приносить дохід, рахунки від провайдера і заходи по відходу від атаки відчутно б'ють «ціль» по кишені. В даний час DoS і DDoS-атаки найбільш популярні, оскільки дозволяють призвести до відмови практично будь-яку систему, не залишаючи юридично значимих доказів [4].

Якщо атака виконується одночасно з великої кількості комп'ютерів, то говорять про DDoS-атаку (від англ. Distributed Denial of Service, розподілена атака типу «відмова в обслуговуванні»). Така атака проводиться в тому випадку, якщо потрібно викликати відмову в обслуговуванні добре захищеної крупної компанії чи державної організації [3].

DDoS – широкомасштабна координована атака на надання послуг системи жертви або мережевих ресурсів, яка побічно запускається через велику кількість комп'ютерних агентів, що потрапили в Інтернет. Перед застосуванням атаки зловмисник приймає велику кількість комп'ютерних машин під його управлінням через Інтернет, і ці комп'ютери є вразливими машинами. Зловмисник використовує недоліки цих комп'ютерів, вставляючи шкідливий код або іншу техніку хакерства, щоб вони стали під його контролем. Ці вразливі або скомпрометовані машини можуть складати сотні або тисячі осіб, і їх зазвичай називають «зомбі». Група зомбі зазвичай

формує «ботнет». Величина атаки залежить від розміру ботнету, для більшого ботнету, атаки є більш серйозними і катастрофічними [8].

Раніше корпоративні комп'ютери часто атакували вірусні програми, які підміняли платіжні доручення, коли бухгалтер намагався провести транзакції, і забирали гроші на підроблені рахунки. Зараз такі програми практично відсутні, але методи шахраїв стали ще більш витонченими. Все частіше стали зустрічатися випадки, коли бухгалтер вставляє спеціальний ключ для доступу до банку, вводить всі паролі, починає проводити транзакцію, а на комп'ютері з'являється картинка, що імітує перезавантаження (програмний код). Насправді за цією картинкою зловмисники використовують вже підготовлену бухгалтером транзакцію для того, щоб перевести гроші на свої рахунки.

Часто злочинці навіть не використовують спеціальні шкідливі програми, обходячись стандартними засобами для віддаленого управління операційною системою, і без всяких картинок підключаються до комп'ютера і переводять гроші. Коли пропажа виявляється, а на комп'ютері немає ніяких вірусів, природно, під підозру відразу потрапляє сам бухгалтер.

Широке поширення отримали програми-вимагачі, які шифрують всі документи на комп'ютері: платіжні доручення, бази даних, звітність, всю документацію, – а для повернення доступу до даних вимагають перерахувати гроші. З корпоративних користувачів вимагають перевести до декількох тисяч доларів або їх еквівалент в біткоїнах.

Фішинг – це спосіб, при якому шахрай може отримати інформацію, не маючи жодного контакту з картою. Вся інформація найчастіше викрадається через Інтернет. У власників можуть вкрати номер карти, термін дії, ПІН-код та CVV/CVC-код. Отримавши всю необхідну інформацію, шахраї з легкістю крадуть гроші з карт. Найбільш поширеним способом фішингу є відправка електронних листів, в яких міститься посилання. Перейшовши по такому посиланню, людина потрапляє на сайт, який нагадує сайт банківської

установи, причому його адресу може відрізнятися від справжнього сайту банку на одну або кілька букв. Неуважний користувач може не помітити підміни і подумати, що це офіційний сайт, надавши йому всю конфіденційну інформацію з карти [20].

Представимо наочно масштаби кібернетичних загроз у банківській системі світу ґрунтуючись на Звіті про тенденції «Фінансові кібернетичні загрози першого кварталу 2017 року», який був розроблений Лабораторією Касперського та компанією Telefónica [13]. В звіті використовуються дані Kaspersky Security Network (KSN) – глобального сервісу оперативної реакції на загрози. Коли програма виявляє підозрілі або неперевірені дані на комп'ютері учасника KSN – ці дані автоматично відправляються в вірусну лабораторію Kaspersky. Часовий інтервал для проведеного аналізу містить дані, отримані в період з 1 січня 2017 року по 31 березня 2017 року.

Станом на кінець першого кварталу 2017 року найбільшої шкоди від фішингових атак зазнають банки – 51,70% (рисунок 1.1).

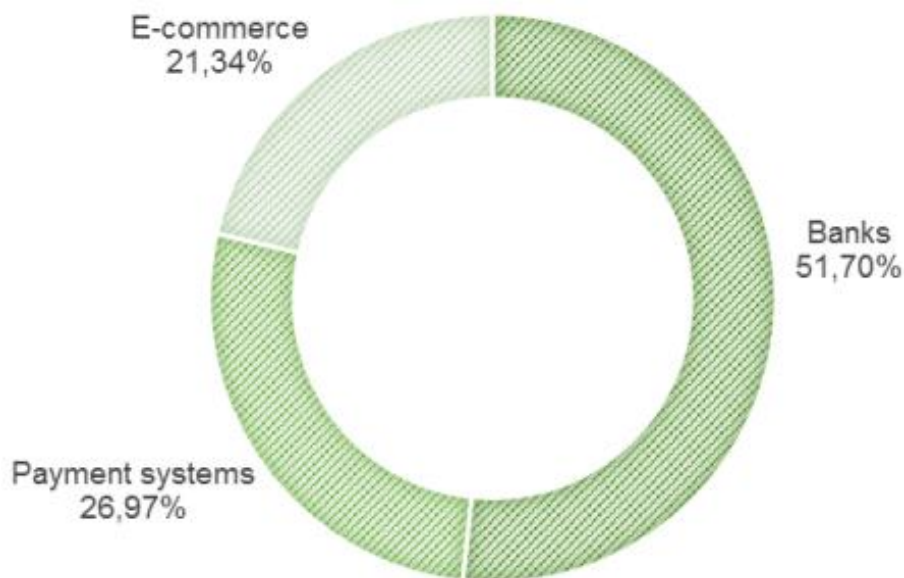


Рисунок 1.1 – Цільовий розподіл фішингу у фінансовому секторі

Так, кількість фішингових атак у фінансовій сфері, зареєстрованих Лабораторією Касперського, скоротилася на 7,1% порівняно з попереднім

кварталом; зменшення частки нападів на банківські установи склало -2,53%. Як і в попередньому періоді, найбільше від фішингу страждають користувачі в Китаї та Бразилії. За ними слідують жителі Макао, Російської Федерації та Австралії.

Наведена нижче карта показує країни з найбільшим відсотком кількості користувачів, які стали жертвами фішингових атак (відношення атакованих користувачів до загальної кількості користувачів KSN у країні, на пристроях із включеними компонентами захисту від фішингу) (рис. 1.2).

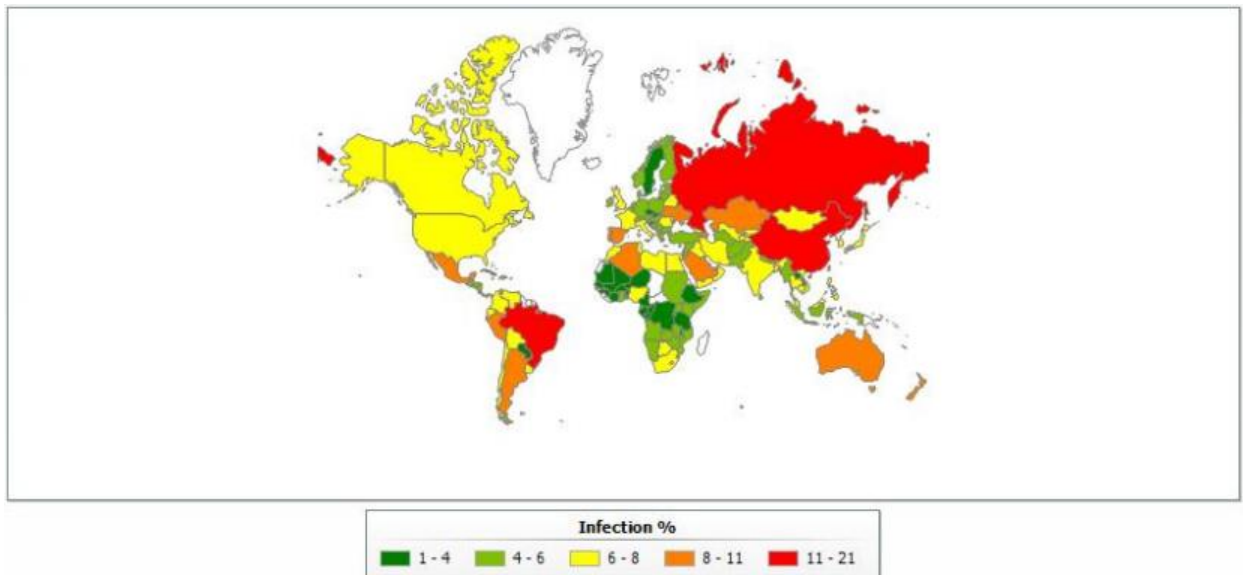


Рисунок 1.2 – Географічне розповсюдження фішингу – перший квартал 2017 року

Наведений нижче графік показує динаміку частки унікальних користувачів по всьому світу, які стали жертвами фішингових атак у першому кварталі 2017 року (рис. 1.3). Як і в попередніх кварталах, графік показує коливання, які відповідають окремим фішинговим компаніям.

Країни з найвищим відсотком нападу на користувачів – Китай (20,87%) та Бразилія (19,16%). За ними слідують Макао (11,94%), Російська Федерація (11,29%) та Австралія (10,73%).

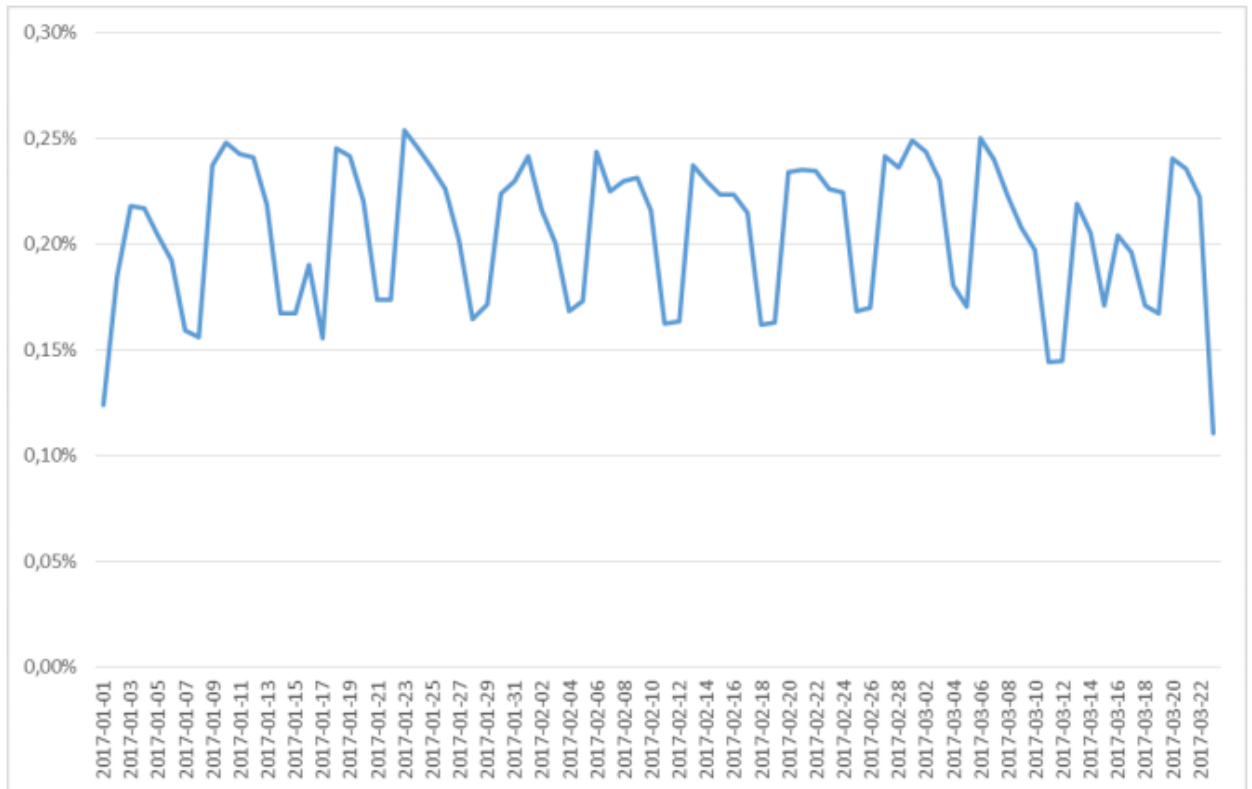


Рисунок 1.3 – Статистика фішингових атак – перший квартал 2017 року

Наступний графік показує відсоток користувачів, які стали жертвами фішингових атак у країнах з найбільшим відсотком атакованих користувачів (рисунок 1.4).

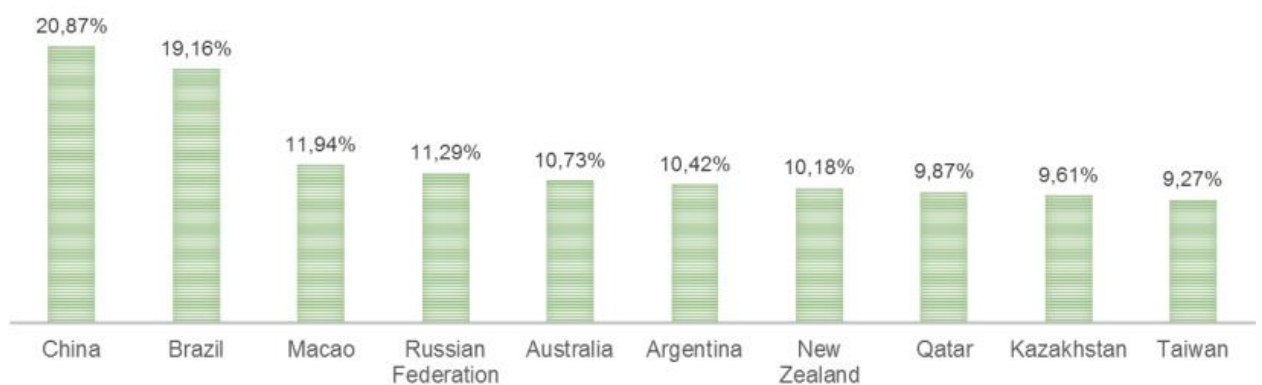


Рисунок 1.4 – Країни з найвищим рівнем жертв від фішингу

Найбільш поширеною мобільною кіберзагрозою є банківські трояни, оскільки в більшості володарів смартфонів є в наявності і банківська карта. А

оскільки банки використовують мобільні номери для авторизації (наприклад, відправляють SMS з одноразовими паролями для підтвердження операцій), в шахраїв виникає спокуса цей канал комунікації перехопити і здійснювати перекази і платежі з чужого банківського рахунку.

Основних методів роботи банківських троянців три:

- вони можуть приховувати від користувача банківські SMS з паролями і тут же перенаправляти їх зловмисникові, який скористається ними, щоб перевести гроші на свій рахунок;
- банківські трояни можуть діяти в автоматичному режимі, час від часу відправляючи відносно невеликі суми на рахунок злочинців;
- зловредів відразу маскують під мобільні додатки банків і, отримавши доступ до реквізитів для входу в мобільний інтернет-банк, роблять все те ж саме [18].

За даними Лабораторії Касперського Banker.AndroidOS.Asacub.ar став найпопулярнішим троянським оператором мобільного зв'язку в третьому кварталі 2017 року, замінивши довгострокового лідера Trojan-Banker.AndroidOS.Svpeng.q. Ці мобільні банківські троянські програми використовують фішингові вікна, щоб викрасти дані кредитної картки, логіни та паролі для онлайн-ових банківських рахунків. Крім того, вони викрадають гроші за допомогою послуг SMS, включаючи мобільний банкінг.

Географія загроз мобільного банкінгу у 3-му кварталі 2017 року (відсоток від усіх атакованих користувачів) зображена на рисунку 1.5.

Частка атакованих користувачів виражена відсотком унікальних користувачів у кожній країні, що зазнали атаки мобільних банківських троянських програм відносно всіх користувачів мобільного продукту безпеки компанії Лабораторії Касперського у країн [5].

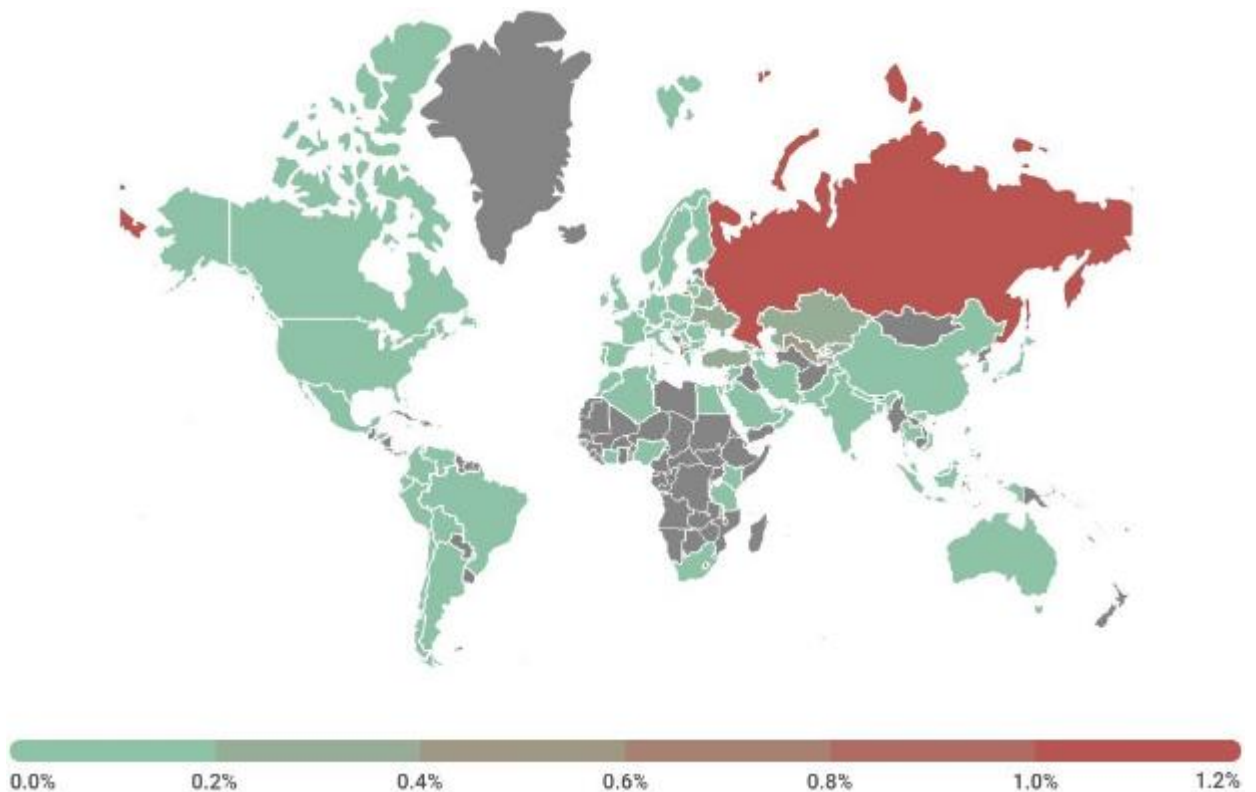


Рисунок 1.5 – Географія загроз мобільного банкінгу у 3-му кварталі 2017

Топ-10 країн, яких атакували мобільні банківські троянські програми (за рейтингом за наслідками атаки користувачів) представлені в таблиці 1.1.

Таблиця 1.1 – Топ-10 країн, атакованих банківськими троянами

№	Країна	Частка атакованих користувачів, %
1	Росія	1,20
2	Узбекистан	0,40
3	Казахстан	0,36
4	Таджикистан	0,35
5	Туреччина	0,34
6	Молдова	0,31
7	Україна	0,29
8	Киргизстан	0,27
9	Білорусь	0,26
10	Латвія	0,23

Розглянемо основні категорії «фізичних» атак (пошкодження або відкриття пристрою, підключення зовнішніх пристроїв), які є традиційними.

Скіммінг – встановлення спеціальних технічних засобів, причому не обов’язково в картоприймач, для розкрадання даних, записаних на магнітну стрічку платіжної картки. PIN-код, як правило, викрадають за допомогою окремого технічного пристрою – відеокамери або фальшивої накладки на PIN-пад.

У ряді випадків відзначено використання нового виду скімінгового обладнання – так званого перископного.

Шиммінг – встановлення в картоприймач спеціальних технічних засобів, призначених для розкрадання даних з EMV-чіпа карти. Таким чином викрадається наступна інформація: історія платежів, інформація, що міститься на Track 2 карти, термін дії.

Black Box – встановлення або підключення технічного пристрою, що взаємодіє з компонентами банкомату (найчастіше з дозатором) і віддає останньому команду для видачі грошових коштів.

Атаки на безконтактні карти (NFC) – створення дублікатів платіжних карт, технічне розкрадання безконтактним методом ряду важливих даних, включаючи тип використовуваного платіжного додатка, термін дії карти, ім’я власника картки, PAN (Primary Account Number) карти та ін.

Підміна процесингу – в цьому випадку банкомат відключається від процесингу кредитної організації і підключається до пристрою, що імітує його. Передові пристрої можуть імітувати нормальний стан банкомату (обслуговування клієнтів) для моніторингу ПЗ. Сутність атаки полягає в передачі банкомату підроблених команд про видачу грошових коштів без порушення загальної логіки роботи банкомату і модифікації його компонентів, як апаратних, так і програмних.

Transaction Reversal Fraud (TRF) – отримання готівкових коштів з одночасним впливом на роботу банкомату і процесингового центру, в результаті чого відсутня коректне завершення операції з видачі готівки й не змінюється баланс по карті (маніпулювання картковим рахунком) [23].

Постійний розвиток комп'ютерних технологій, без яких не може обійтись жоден банк, призводить до появи все більшої кількості нових кіберзагроз в банківській сфері. У зв'язку з чим постає питання стосовно необхідності виявлення та попередження цих загроз.

1.2 Аналіз існуючих підходів до виявлення кіберзагроз у банках

Кібербезпека стала дуже актуальною і затребуваною в Україні. Така сучасна реальність, що виникла в результаті масштабних хакерських атак.

2017 рік став переломним і дуже складним для банківської системи України. Навесні багато банків виявилися зараженими глобальним вірусом шифрувальником WannaCry, який вражає комп'ютери з операційною системою Microsoft Windows. Ця шкідлива програма масово виводила з ладу робочі місця касирів, менеджерів, операціоністів та інших банківських службовців.

Здається, що після цього випадку всі, хто постраждав, повинні були прийняти серйозні заходи для підвищення безпеки банківських систем і мереж. Але реальність показала наступне: 27 червня в 10.00 вірус Petya.A вразив 70% банківської інфраструктури України.

Це призвело до того, що деякі системоутворюючі банки кілька годин не могли провести платежі. Цього разу постраждали не тільки робочі місця банківських службовців, але й, що дуже небезпечно, сервери та бази даних. Через місяць після закінчення атаки вірусу Petya.A деякі банки ще продовжували відновлювати порушені бізнес-процеси і втрачені дані.

Після цієї масштабної хакерської атаки, яка сильно вразила громадськість, Національний банк України, як регулятор банківської системи, почав розробляти заходи, які дозволили б у майбутньому запобігти подібним ексцесам.

В результаті цієї роботи НБУ було прийнято Постанову №95 «Про затвердження Положення про організацію заходів із забезпечення інформаційної безпеки в банківській системі України» від 28 вересня 2017 року, яку було офіційно опубліковано 04 жовтня 2017 року [35]. Наведений перелік вимог до банківських установ буде перевірятись Національним банком України вже з 01 березня 2018 року.

Ця Постанова розроблена з метою удосконалення вимог до захисту інформації в інформаційних системах банків з урахуванням актуальних кіберзагроз та складається з вимог, які необхідно впровадити для забезпечення інформаційної безпеки і кіберзахисту банків.

Однією з вимог є впровадження банками наступних основних технічних систем:

- виявлення атак;
- моніторинг події управління інцидентами;
- контроль доступу до мережі;
- захист електронної пошти;
- запобігання атак, спрямованих на відмову в обслуговуванні;
- антивірусний захист;
- двофакторна аутентифікація.

Наступні вимоги мають на увазі виконання організаційних змін або створення нових бізнес-процесів. Так, наприклад, банк зобов'язаний сформувати керівний орган з інформаційної безпеки, в який необхідно включити вище керівництво рівня правління. Також кожен банк, незалежно від його розміру, повинен мати підрозділ, що відповідає за інформаційну безпеку, і в його складі має бути не менше двох осіб.

Кожен український банк зобов'язаний регулярно проводити тестування на проникнення в критичні системи банку.

Остання група вимог – необхідно, щоб в банках були розроблені нормативні документи, що встановлюють правила для роботи персоналу.

Аспекти, які повинні бути прописані в цих документах:

- використання змінних носіїв інформації;
- надання, скасування та контроль доступу до банківських інформаційних систем;
- використання криптографії;
- контроль змін на рівні мережі;
- використання електронної пошти.

Для того щоб впровадити всі вимоги Постанови №95, українському банку необхідно мати досвідчених фахівців в різноманітних областях інформаційної безпеки.

У цих фахівців повинен бути великий досвід в побудові системи інформаційної безпеки, в проектуванні процесів, у впровадженні та експлуатації технічних систем захисту, а також в розробці нормативних документів [35].

В останні роки світовий ринок інструментів для боротьби з шахрайством розвивається бурхливими темпами. Очевидно, що сьогодні консервативних методів вже недостатньо в силу постійного вдосконалення шахраями своїх прийомів і схем, тому банкам потрібні якісно інші рішення для проактивного виявлення та запобігання шахрайству.

Багато українських банків все ще покладаються на збільшення штату співробітників служби безпеки або залучення додаткових ресурсів для проведення розслідувань. Безумовно, і НБУ як регулятор, і різні профільні асоціації стимулюють банки вдосконалювати свої методи боротьби з кіберзагрозами, розвиваючи законодавчі норми. Однак цей процес в Україні все ще на стадії формування, а протидія кібернетичним злочинам поки є приватною справою кожної фінансової організації.

Західні банки пішли більш інноваційним шляхом, запровадивши інструменти бізнес-аналітики, завдяки яким можна не просто виявляти шахрайські схеми постфактум, але і вживати превентивних заходів для їх

запобігання. Застосовуючи інструменти поглибленої аналітики, сьогодні банки можуть відслідковувати всю історію взаємодії з кожним клієнтом, що істотно допомагає їм розпізнати потенційне шахрайство. Для цього потрібно враховувати досить великий набір параметрів, проводити аналіз соціального оточення клієнтів і встановлювати найрізноманітніші взаємозв'язки – за адресою проживання, телефонними номерами, кредитною історією, заставному майну і т.д. Іншими словами, технологічно процес протидії кіберзагрозам передбачає регулярний аналіз даних на предмет виявлення як раніше відомих, так і нових випадків подібної поведінки.

Компанія SAS пропонує так званий гібридний підхід для аналізу та виявлення кіберзагроз. Його сутність полягає в комбінуванні різних методів і алгоритмів:

- застосування відомих експертних і статистичних бізнес-правил як прямого, так і нечіткого збігу, які допомагають виявляти невідповідність інформації в різних джерелах даних, невідповідність кодів операцій, дублікати та інші;
- моделі незвичайної «аномальної» поведінки або моделі відхилення від звичайної схеми поведінки клієнта;
- пошук прихованих закономірностей в даних і побудова прогнозних моделей, які дозволяють прогнозувати вже виявлені схеми;
- аналіз соціальних мереж, які допомагають отримати нові дані на підставі аналізу зв'язків, наприклад, зв'язку з відомими випадками кіберзагроз, виявлення маніпулювання, транзакції з підозрілими контрагентами (Social Network Analysis).

Такий підхід дозволяє підвищити кількість знайдених випадків шахрайства і знизити помилку першого роду («хибна тривога») [17].

Актуальність проблем кібербезпеки у банківській сфері не викликає жодних сумнівів. Для повноцінного захисту систем зберігання та обробки даних потрібна не просто установка відповідного програмного забезпечення,

а цілий комплекс програмно-технічних, адміністративно-організаційних та нормативно-правових заходів. Для виявлення кіберзагроз в банку необхідним є опрацювання великої кількості даних щодо минулих транзакцій, які згодом виявилися шахрайськими. Виникає проблема вилучення корисної для користувача інформації з великого обсягу «сирих» даних. Оскільки людина не в змозі самотійно опрацьовувати такі обсяги даних, доцільним є застосування методів інтелектуального аналізу.

1.3 Аналіз методів інтелектуального аналізу

Технології аналізу даних, що базуються на застосуванні класичних статистичних підходів, мають низку недоліків. Відповідні методи ґрунтуються на використанні усереднених показників, на підставі яких важко з'ясувати справжній стан справ у досліджуваній сфері. Методи математичної статистики виявилися корисними насамперед для перевірки заздалегідь сформульованих гіпотез та «грубого» розвідницького аналізу, що становить основу оперативної аналітичної обробки даних (OLAP). Окрім того, стандартні статистичні методи відкидають (нехтують) нетипові спостереження – так звані піки та сплески. Проте окремі нетипові значення можуть становити самотійний інтерес для дослідження, характеризуючи деякі виняткові, але важливі явища [29].

Ці недоліки статистичних методів спонукали до розвитку нових методів дослідження складних систем, які останнім часом все частіше застосовуються для вирішення практичних завдань – методів інтелектуального аналізу даних. Інтелектуальний аналіз даних (далі ІАД) – виявлення прихованих закономірностей або взаємозв'язків між змінними у великих масивах необроблених даних.

Сфера застосування ІАД нічим не обмежена – вона скрізь, де є якісь дані. Але насамперед методи ІАД сьогодні зацікавили комерційні

підприємства, що розгортають свої проекти на основі інформаційних сховищ даних (Data Warehousing). ІАД являють собою велику цінність для керівників і аналітиків у їх повсякденній діяльності. Ділові люди усвідомили, що за допомогою методів ІАД вони можуть одержати відчутні переваги у конкурентній боротьбі. Досвід багатьох підприємств показує, що віддача від використання ІАД може сягати 1000 % [26].

В англomовній літературі замість терміна «інтелектуальний аналіз даних» зазвичай використовується термін Data Mining (дослівний переклад – «видобуток даних»)

Технологію Data Mining достатньо точно визначає Григорій Піатецький - Шапіро (Gregory Piatetsky - Shapiro) – один із засновників цього напрямку: «Data Mining – це процес виявлення в сирих даних раніше невідомих, нетривіальних, практично корисних і доступних інтерпретації знань, необхідних для ухвалення рішень в різних сферах людської діяльності» [16].

В основу інтелектуального аналізу покладена концепція шаблонів (паттернів), що відбивають фрагменти багатоаспектних взаємин у даних. Ці шаблони являють собою закономірності, властиві підвибіркам даних, які можуть бути компактно виражені у зрозумілій людині формі. Пошук шаблонів проводиться методами, не обмеженими апріорними припущеннями про структуру вибірки, та видами розподілів значень аналізованих показників. Важливе положення інтелектуального аналізу – нетривіальність розшукуваних шаблонів. Це означає, що знайдені шаблони повинні відбивати неочевидні, несподівані регулярності в даних, що становлять так звані приховані знання. До суспільства прийшло розуміння, що сирі дані містять глибинний шар знань, за грамотного «розкопування» якого можуть бути виявлені справжні «самородки» [26].

Розглянемо основні задачі, які вирішуються методами Data Mining:

- класифікація – віднесення об'єктів (спостережень, подій) до одного з заздалегідь відомих класів;
- регресія (в тому числі задачі прогнозування) – встановлення залежності безперервних вихідних від вхідних змінних;
- кластеризація – угруповання об'єктів (спостережень, подій) на основі даних (властивостей), що описують сутність цих об'єктів, у кластери. Об'єкти всередині кластера повинні бути «схожими» один на одного і відрізнятися від об'єктів, що ввійшли в інші кластери. Чим більше схожі об'єкти всередині кластера і чим більше відмінностей між кластерами, тим точніша кластеризація;
- асоціація – виявлення закономірностей між пов'язаними подіями. Прикладом такої закономірності служить правило, яке вказує, що з події X слід подія Y. Такі правила називаються асоціативними. Вперше ця задача була запропонована для знаходження типових шаблонів покупок, що здійснюються в супермаркетах, тому іноді її ще називають аналізом ринкової кошика (market basket analysis);
- послідовні шаблони – встановлення закономірностей між пов'язаними в часі подіями, тобто виявлення залежності, що якщо відбудеться подія X, то через заданий час відбудеться подія Y;
- аналіз відхилень – виявлення найбільш нехарактерних шаблонів [28].

Банківська сфера в усьому світі зазнала величезних змін в шляхах ведення бізнесу. Використання «електронного банківського обслуговування» призводить до простішого захоплення транзакційних даних та одночасного збільшення обсягу таких даних. Люди самотійно не в змозі аналізувати цю величезну кількість необоротних даних та ефективно перетворити їх в корисні для організації знання [9].

Інтелектуальний аналіз даних може допомогти у вирішенні проблем бізнесу шляхом пошуку моделей, асоціацій, кореляцій, які приховані в діловій інформації, що зберігається в базах даних банків.

Методи інтелектуального аналізу можна розподілити на технологічні, статистичні та кібернетичні. Розглянемо більш детально ці методи в таблиці 1.2.

Таблиця 1.2 – Методи інтелектуального аналізу [21]

Група методів	Опис
Технологічні методи	а) безпосереднє використання даних, або збереження даних. Методи цієї групи: кластерний аналіз, метод найближчого сусіда; б) виявлення і використання формалізованих закономірностей, або дистиляція шаблонів - логічні методи, методи візуалізації, методи крос-табуляції, методи, що засновані на рівняннях.
Статистичні методи	а) описовий аналіз і опис вихідних даних; б) аналіз зв'язків (кореляційний і регресійний аналіз, факторний аналіз, дисперсійний аналіз); в) багатовимірний статистичний аналіз (компонентний аналіз, дискримінантний аналіз, багатовимірний регресійний аналіз, канонічні кореляції і ін.); г) аналіз тимчасових рядів (динамічні моделі і прогнозування).
Кібернетичні методи	а) штучні нейронні мережі (розпізнавання, кластеризація, прогноз); б) еволюційне програмування (в т.ч. алгоритми методу групового обліку аргументів); в) генетичні алгоритми (оптимізація); г) асоціативний алгоритм; д) нечітка логіка; е) дерева рішень; ж) системи обробки експертних знань.

Традиційні методи аналізу даних в основному орієнтовані на перевірку наперед сформульованих гіпотез (статистичні методи) і на «грубий розвідувальний аналіз», що становить основу оперативної аналітичної обробки даних (Online Analytical Processing, OLAP), тоді як одне з основних положень Data Mining – пошук неочевидних закономірностей. Інструменти Data Mining можуть знаходити такі закономірності самостійно і також самостійно будувати гіпотези про взаємозв'язки. Оскільки саме

формулювання гіпотези щодо залежності є найскладнішою задачею, перевага Data Mining в порівнянні з іншими методами аналізу є очевидною.

Більшість статистичних методів для виявлення взаємозв'язків в даних використовує концепцію усереднювання по вибірці, що приводить до операцій над неіснуючими величинами, тоді як Data Mining оперує реальними значеннями.

OLAP більше підходить для розуміння ретроспективних даних, Data Mining спирається на ретроспективні дані для отримання відповідей на питання про майбутнє.

Методи Data Mining можна поділити на дві групи:

- Supervised Learning (навчання з вчителем);
- Unsupervised Learning (навчання без вчителя).

У першому випадку завдання аналізу даних, наприклад класифікація, здійснюється в кілька етапів. Це один із способів машинного навчання, в ході якого випробувана система примусово навчається за допомогою прикладів. Між входами і еталонними виходами може існувати деяка залежність, але вона не відома. Відома тільки кінцева сукупність прецедентів, звана навчальною вибіркою. На основі цих даних потрібно відновити залежність, тобто побудувати алгоритм, здатний для будь-якого об'єкта видати досить точну відповідь. Для вимірювання точності відповідей, так само як і в навчанні на прикладах, може вводитися функціонал якості. Спочатку за допомогою будь-якого алгоритму Data Mining будується модель аналізованих даних – класифікатор. Потім класифікатор піддається «навчанню». Іншими словами, перевіряється якість його роботи і, якщо вона незадовільна, відбувається «додаткове навчання» класифікатора. Так триває до тих пір, поки не досягнемо необхідного рівня якості або не переконаємося, що обраний алгоритм не працює коректно з даними або ж самі дані не мають структури, яку можна виявити.

Unsupervised Machine Learning – один із способів машинного навчання. З його допомогою випробувана система спонтанно навчається виконувати поставлене завдання без втручання з боку експериментатора. Як правило, це придатне тільки для завдань, в яких відомі описи безлічі об'єктів (навчальної вибірки) і потрібно виявити внутрішні взаємозв'язки, залежності, закономірності, що існують між об'єктами. Наприклад, закономірності в покупках, скоєних клієнтами великого магазину. Очевидно, що якщо ці закономірності існують, то модель повинна їх представити і недоречно говорити про її навчання. Звідси і назва «навчання без вчителя».

Навчання без вчителя часто протиставляється навчанню з вчителем, коли для кожного навчального об'єкта примусово задається «правильна відповідь» і потрібно знайти залежність між стимулами і реакціями системи [30].

Найбільш перспективним напрямком інтелектуального аналізу є кібернетичні методи, які представляють собою множину підходів, об'єднаних ідеєю комп'ютерної математики та використання теорії штучного інтелекту [28]. До основних методів цієї групи належать нейронні мережі, дерева рішень, кластерний та асоціативний аналізи, побудова нелінійних регресій та інші. Використання інтелектуального аналізу допомагає у побудові концептуальної моделі, головною метою якої є виявлення кібернетичних загроз в банківських установах для майбутнього попередження цих загроз в разі їх виникнення.

1.4 Побудова концептуальної моделі виявлення ознак кіберзагроз

Головна мета побудови концептуальної моделі – виявлення ознак кібернетичних загроз в банківських установах для майбутнього попередження цих загроз в разі їх виникнення.

Побудова моделі інтелектуального аналізу даних є складовою частиною масштабнішого процесу, який включає всі етапи, починаючи з визначення базової проблеми, яку модель вирішуватиме, до розгортання моделі в робочому середовищі. Даний процес може бути заданий за допомогою наступних шести базових кроків:

- постановка задачі;
- підготовка даних;
- перегляд даних;
- побудова моделей;
- дослідження, перевірка, прогнозування за допомогою моделей;
- розгортання і оновлення моделей [22].

Розглянемо алгоритм моделювання ймовірності виникнення ознак кібернетичної загрози в банківських транзакціях:

- визначення вхідних даних моделі виявлення ймовірності виникнення ознак кіберзагроз, в якості яких виступатиме інформація по транзакціям користувачів мобільного та інтернет-банкінгу;
- дослідження вхідних даних для виявлення тенденцій, взаємозв'язків та розподілу, аналіз отриманих результатів, за необхідністю проведення коригування змінних шляхом їх логарифмування;
- розробка математичної моделі у вигляді формалізованого представлення математичних залежностей, які будуть описувати вхідні дані, та на основі цього, моделювати результуючу величину;
- програмна реалізація моделі, яка включає вибір програмних засобів, за допомогою яких буде реалізовано модель;
- отримання вихідних даних – ймовірності виникнення ознак кіберзагроз;
- перевірка якості й адекватності моделі – це відповідність моделі фактам і тенденціям реального економічного буття;

- інтерпретація результатів моделювання – отримані числові результати інтерпретуються з точки зору їх економічного змісту;
- формулювання висновків – результати моделювання аналізуються та узагальнюються тенденції [11].

Побудова даної моделі передбачає використання наступних даних банку, що містять інформацію про проведені транзакції користувачами інтернет-банкінгу або мобільного банкінгу: їх суми, частоту, географічне положення, перевищення встановлених лімітів, обнуління рахунків та іншу.

Вибір цих факторів обумовлений тим, що користувачі мобільного та інтернет-банкінгу є однією із слабких ланок в системі банківської безпеки. Це пов'язано з тим, що банк не в змозі контролювати, хто є користувачем та де він користується пристроєм. Частіше за все такі операції можуть містити ознаки кіберзагрози, тобто піддатися під різновид соціальної інженерії.

Для побудови моделі висунуто ряд гіпотез стосовно вірогідності виникнення ознак кіберзагроз під час проведення транзакцій користувачами мобільного та інтернет-банкінгу. Показники, що можуть вказувати на можливе виникнення кіберзагрози в процесі виконання банківської операції:

1. Транзакція має ознаки кіберзагрози, якщо її ініційовано на території іншої країни.

В більшості банків прийнята практика необхідності повідомлення банку клієнтом про виїзд за кордон та зазначення країн, які будуть відвідані. В іншому випадку служба безпеки банку може заблокувати карту, якщо по ній будуть ініціюватися транзакції в іншій країні. Це пов'язано з тим, що хакери, зламуючи доступ до мобільного або інтернет-банкінгу та привласнюючи чужі кошти, застосовують спеціальні програми для шифрування їх місцеположення.

2. Від типу пристрою, з якого виконувалась транзакція залежить ймовірність виникнення кіберзагрози.

Існують різні способи злому мобільних пристроїв та комп'ютерів, завдяки яким зловмисники з легкістю отримують доступ до мобільного та інтернет-банкінгу користувачів банківських послуг.

3. На ймовірність виникнення ознак кіберзагрози впливає тип проведеної транзакції.

Широке розповсюдження типів банківських транзакцій сприяє впровадженню нових заходів з боку зловмисників, направлених на заволодіння чужими коштами та порушення безпеки інформації в банку.

4. Онуління рахунків клієнтів банку вказує на ймовірні ознаки кіберзагроз.

В наш час досить розповсюдженими є безготівкові розрахунки, коли платежі відбуваються без використання готівкових коштів. Тому, в більшості випадків на банківському рахунку людини завжди присутня певна сума коштів, під час транзакції зі зняття всієї суми ймовірно має місце ознака порушення користування рахунком, можливе несанкціоноване зняття коштів.

В процесі застосування моделі відбувається перевірка наявності потенційної ознаки кіберзагрози в транзакції, що виконується. Якщо є підозра на кібернетичну загрозу, банк має повідомити про це клієнта банку, тимчасово призупинивши транзакцію та затребувавши від користувача інтернет-банкінгу підтвердження операції, шляхом введення захисного коду, отриманого в СМС чи за телефонним дзвінком від працівника банку.

Виходячи з обраних вхідних даних та методів інтелектуального аналізу було розроблено концептуальну модель виявлення ознак кіберзагроз в транзакціях користувачів мобільного та інтернет-банкінгу, зображену на рисунку 1.6.

Таким чином, використання моделі дасть змогу попередити типові кіберзагрози, метою яких є порушення прав клієнтів та цілісності, безпеки інформації.



Рисунок 1.6 – Концептуальна модель виявлення ознак кіберзагроз в банківських транзакціях

Проте дана модель потребує постійного оновлення та удосконалення у зв'язку з появою нових загроз для користувачів мобільного та інтернет-банкінгу. Необхідно доповнювати вибірку даних актуальною інформацією про виконані користувачами транзакції.

Для більш ефективної роботи моделі необхідно інтегрувати її з діючою банківською системою.

РОЗДІЛ 2 МАТЕМАТИЧНІ МОДЕЛІ ВИЯВЛЕННЯ ОЗНАК КІБЕРЗАГРОЗ У БАНКАХ

2.1 Первинний аналіз вхідних та вихідних даних

В процесі підготовки до побудови моделі виявлення ознак кіберзагроз у банку в якості вихідних даних було використано інформацію, що міститься у базі даних мобільного та інтернет-банкінгу банку «Х». Оскільки дана інформація є комерційною таємницею, то розголошення назви банківської установи не є можливим. Інформація містить 8 вхідних змінних, включаючи цільову змінну. Назви, зміст, ролі та типи змінних представимо в таблиці 2.1.

Таблиця 2.1 – Опис вхідних змінних

Ім'я змінної	Економічний зміст	Роль	Тип	Допустимі значення
isfraud (Y)	Випадки виникнення кіберзагроз	цільова	binary	1 – виявлено ознаки кіберзагроз; 0 – ознак кіберзагроз не виявлено.
amount (X_1)	Загальна сума, що проходила в транзакціях	вхідна	interval	≥ 0
devicetype (X_2)	Тип пристрою, з якого виконувалась транзакція	вхідна	nominal	М – мобільний банкінг; І – інтернет банкінг.
factlocation (X_3)	Ініційоване місцеположення пристрою, з якого проводилась транзакція	вхідна	nominal	UA – Україна; Other – інша країна.
location (X_4)	Місцеположення, вказане при реєстрації клієнта банкінгу	вхідна	nominal	UA – Україна.
newbalance (X_5)	Баланс клієнта після проведення транзакції	вхідна	interval	≥ 0
oldbalance (X_6)	Баланс клієнта до проведення транзакції	вхідна	interval	≥ 0
type (X_7)	Тип виконаної транзакції	вхідна	nominal	CASH_IN – поповнення коштів; CASH_OUT – зняття коштів; DEBIT – списання коштів з рахунку; PAYMENT – проведення оплати; TRANSFER – переведення коштів.

Вибірка даних склала 200000 спостережень, взятих на прикладі інформації за транзакціями користувачав мобільного та інтернет-банкінгу банку «А».

Змінна Y надає дані про те, чи мають місце в банківській транзакції ознаки кіберзагроз, виходячи з інформації за відповідною транзакцією.

Змінна X_1 представлена загальною сумою, що використана певним користувачем банку під час проведення різноманітних транзакцій.

Змінна X_2 вказує на тип пристрою, з якого було проведено транзакцію: мобільний банкінг – мобільний телефон; інтернет-банкінг – комп'ютер.

Змінна X_3 відображає ініційоване місцеположення пристрою, з якого проведено транзакцію: Україна або інша країна.

Змінна X_4 показує, яка країна була вказана користувачем мобільного або інтернет-банкінгу при реєстрації.

Змінна X_5 містить суму, що знаходиться на балансі клієнта після проведення транзакції.

Змінна X_6 містить суму, що знаходилась на балансі клієнта до проведення транзакції.

Змінна X_7 надає інформацію про тип транзакції, яку було проведено користувачем мобільного або інтернет-банкінгу.

Проаналізуємо вхідні дані для виявлення певних закономірностей і тенденцій. На рисунку 2.1 зобразимо кругову діаграму розподілу транзакцій за ймовірністю виникнення ознак кіберзагроз. Серед набору вхідних даних про банківські операції, 22% транзакцій мають ознаки кібернетичних загроз, а у 78% – ознак кіберзагроз не виявлено. Тобто майже 1/5 всієї вибірки має ознаки кібернетичних загроз.

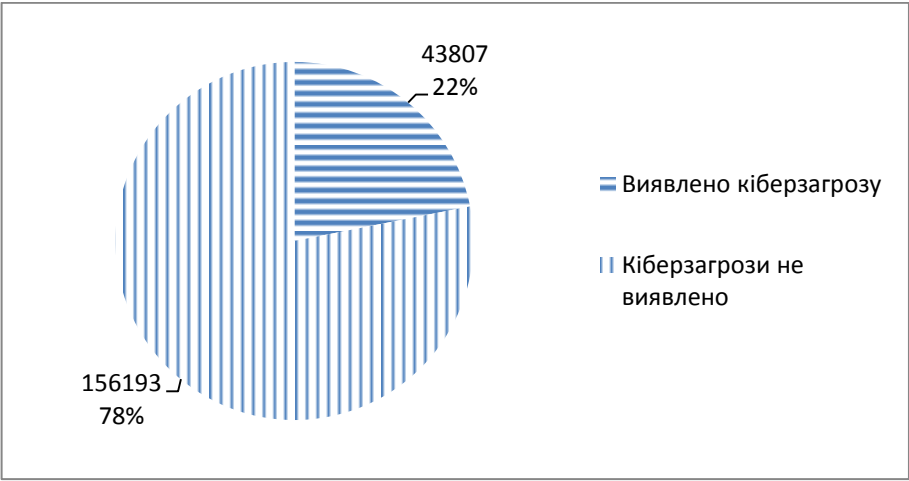


Рисунок 2.1 – Розподіл транзакцій за ймовірністю виникнення ознак кіберзагроз

На рисунку 2.2 представимо розподіл банківських транзакцій за їх типами. Найбільшу долю серед проведених транзакцій становлять проведення оплати (37%), зняття коштів (33%) та поповнення коштів (21%). Незначна частка транзакцій приходить на переведення коштів (8%) та списання коштів (1%).

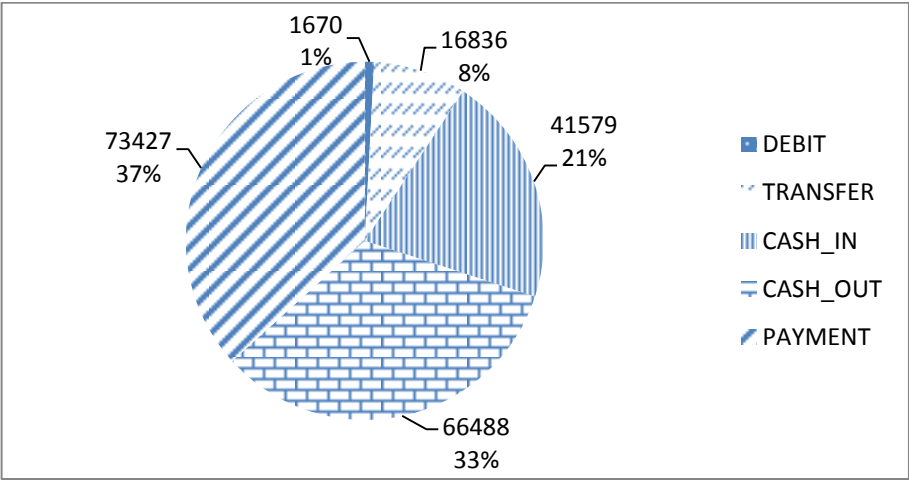


Рисунок 2.2 – Розподіл транзакцій за їх типами

На рисунку 2.3 зобразимо розподіл банківських транзакцій за типами пристроїв, з яких вони виконувались. Розподіл пристроїв мобільного (51%) та інтернет-банкінгу (49%) майже однаковий.

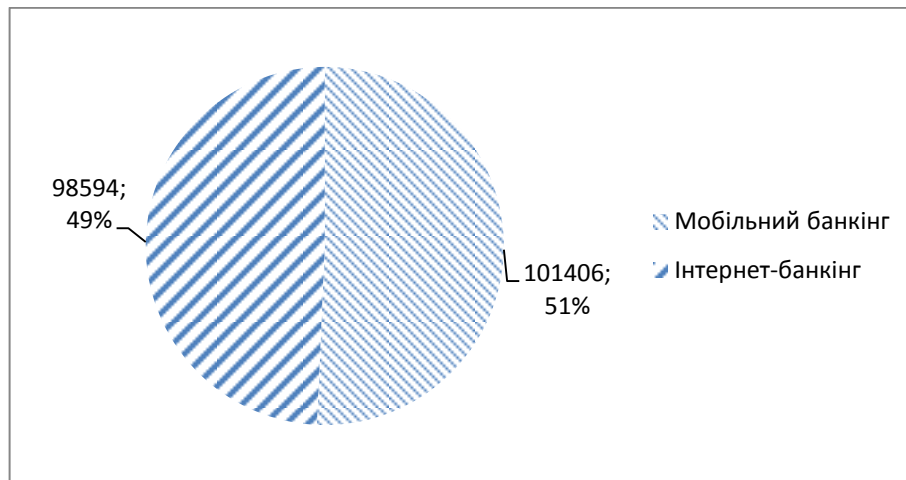


Рисунок 2.3 – Розподіл транзакцій за типами пристроїв, з яких вони виконувались

На рисунку 2.4 представимо розподіл банківських транзакцій за місцеположенням пристрою, з якого проводилась транзакція. В більшості виконаних транзакцій (78%) місцеположення пристрою визначалось як Україна, 22% транзакцій було зафіксовано в інших країнах.

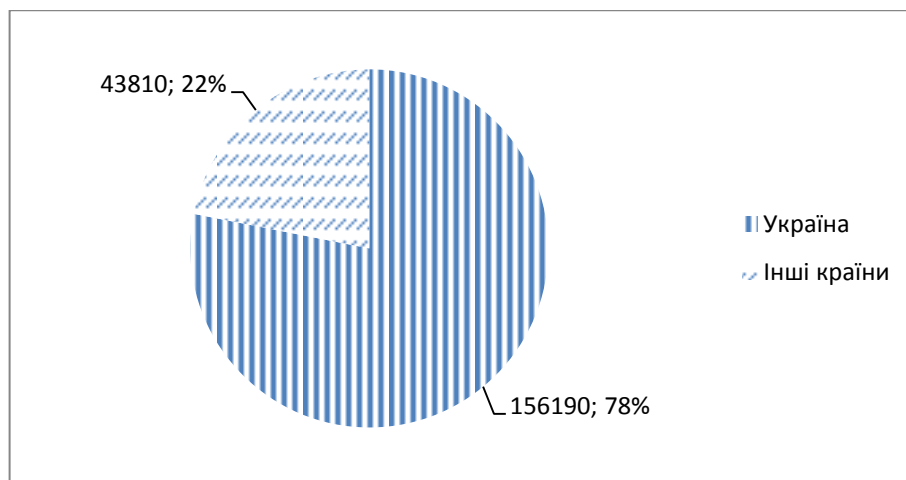


Рисунок 2.4 – Розподіл транзакцій за місцеположенням пристрою, з якого проводилась транзакція

Таким чином, було обрано вхідні змінні для подальшого їх застосування з метою побудови моделей виявлення кіберзагроз в банках методами інтелектуального аналізу.

2.2 Кластерний аналіз як інструмент дослідження первинних даних

Задача кластеризації подібна до задачі класифікації, є її логічним продовженням, але її відмінність в тому, що класи набору даних, що вивчається заздалегідь не визначені.

Мета кластеризації – пошук існуючих структур. Кластеризація є описовою процедурою, вона не робить жодних стратегічних висновків, проте дає можливість провести розвідчий аналіз та вивчити структуру даних.

Нехай X – множина об'єктів, Y – множина номерів (імен, міток) кластерів. Задана функція відстані між об'єктами $\rho(x, x') \in \mathbb{R}$ кінцева навчальна вибірка об'єктів $X^m = \{x_1, x_2, \dots, x_m\} \in X$. Потрібно розбити вибірку на непересічні підмножини, які називаються кластерами, так, щоб кожен кластер складався з об'єктів, близьких за метрикою ρ , а об'єкти різних кластерів істотно відрізнялися. При цьому кожному об'єкту $x_i \in X^m$ приписується номер кластера y_i .

Алгоритм кластеризації – це функція $\alpha: X \rightarrow Y$, яка будь-якому об'єкту $x \in X$ ставить у відповідність номер кластера $y \in Y$. Множина Y в деяких випадках відома заздалегідь, однак частіше ставиться завдання визначити оптимальне число кластерів, з точки зору того чи іншого критерію якості кластеризації.

Кластер можна охарактеризувати як групу об'єктів, що мають спільні властивості. Характеристиками кластера можна назвати дві ознаки:

- внутрішня однорідність;
- зовнішня ізоляваність.

Існує велика кількість підходів до кластеризації:

- алгоритми, засновані на поділі даних (Partitioning algorithms), в тому числі ітеративні: поділ об'єктів на k кластерів; ітеративний перерозподіл об'єктів для поліпшення кластеризації;
- ієрархічні алгоритми (Hierarchy algorithms);

- методи, засновані на концентрації об'єктів (Density-based methods);
- ґрид-методи (Grid-based methods);
- модельні методи (Model-based).

Слід зазначити, що в результаті застосування різних методів кластерного аналізу можуть бути отримані кластери різної форми. В результаті застосування різних методів кластеризації можуть бути отримані неоднакові результати, це є особливістю роботи того чи іншого алгоритму. Однак створення подібних кластерів різними методами вказує на правильність кластеризації.

Задачі кластерного аналізу можна об'єднати в наступні групи:

- розробка типології або класифікації;
- дослідження корисних концептуальних схем групування об'єктів;
- представлення гіпотез на основі дослідження даних;
- перевірка гіпотез або досліджень для визначення, чи дійсно типи (групи), виділені тим чи іншим способом, присутні в наявних даних.

Як правило, при практичному використанні кластерного аналізу одночасно вирішується кілька із зазначених задач [39].

Досліджуючи один або більше атрибутів або класів, можна згрупувати окремі елементи даних разом, отримуючи структурований вивід. На простому рівні при кластеризації використовується один або декілька атрибутів в якості основи для визначення кластера подібних результатів. Кластеризація корисна при визначенні різної інформації, тому що вона корелюється з іншими прикладами так, що можна побачити, як подібність і діапазони узгоджуються між собою. Метод кластеризації працює в обидві сторони. Можна припустити, що в певній точці мається кластер, а потім використовувати свої критерії ідентифікації, щоб перевірити це [14].

У непараметричному випадку ми не маємо інформації про загальний вигляд функцій $f_j(X, \Theta_j)$. Ми можемо мати лише окремі загальні відомості про них: компактність або обмеженість діапазонів змінювання компонент

класифікованих багатовимірних спостережень, неперервність або гладкість відповідних законів розподілу ймовірностей тощо. Вихідні дані зазвичай подають у вигляді матриці спостережень, яка містить значення всіх ознак для кожного із досліджуваних об'єктів, або матриці подібності, що містить попарні відстані між класифікованими спостереженнями.

Бажано, щоб компоненти вектора X відповідали одному й тому самому типу даних. Для цього зазвичай використовують перехід від кількісних ознак до порядкових та від порядкових до номінальних. Але слід урахувати, що при цьому втрачається частина корисної інформації.

Для формалізації задачі класифікації кожний об'єкт зручно інтерпретувати як точку в багатовимірному просторі ознак. Геометрична близькість точок у такому просторі відповідає близькості досліджуваних об'єктів з погляду досліджуваних властивостей.

Класичними непараметричними методами класифікації без навчання є методи кластерного аналізу (таксономії). За їх допомогою вирішують проблему такого розбиття (класифікації, кластеризації) множини об'єктів, за якого всі об'єкти, що належать до одного класу, були б більш подібними один до одного, ніж до об'єктів інших класів. З формальної точки зору, основне завдання методів кластерного аналізу можна сформулювати, як визначення класів еквівалентності й рознесення за ними досліджуваних об'єктів. Під класом, як правило, розуміють генеральну сукупність, що описується одномодальною функцією щільності ймовірності $f(X)$ або, у випадку дискретних ознак, – одномодальним полігоном імовірностей. Номери класів не мають змістового навантаження й використовуються лише для того, щоб відрізнити їх один від одного.

Для формування кластерів застосовують міри подібності та відмінності даних, які можуть бути поділені на три основних види:

- міри подібності (відмінності) типу «відстань» (при їх застосуванні об'єкти вважають тим більш подібними один до одного, чим меншою є відстань між ними);
- міри подібності типу «зв'язок» (у цьому випадку об'єкти вважають тим більш подібними, чим сильнішим є зв'язок між ними);
- інформаційна статистика [15].

Як і будь-які інші методи, методи кластерного аналізу мають певні слабкі сторони, тобто деякі складності, проблеми та обмеження. При проведенні кластерного аналізу слід враховувати, що результати кластеризації залежать від критеріїв розбиття сукупності вихідних даних. При зниженні розмірності даних можуть виникнути певні спотворення, за рахунок узагальнень можуть загубитися деякі характеристики об'єктів.

Існує ряд складнощів при проведенні кластеризації:

1. Складність вибору характеристик, на основі яких проводиться кластеризація. Необдуманий вибір призводить до неадекватного розбиття на кластери і, як наслідок, – до невірної рішення задачі;
2. Складність вибору методу кластеризації. Цей вибір вимагає хорошого знання методів і передумов їх використання. Щоб перевірити ефективність конкретного методу в певній предметній області, доцільно застосувати таку процедуру: розглядають кілька апріорі різних між собою груп і перемішують їх представників між собою випадковим чином. Далі проводиться кластеризація для відновлення вихідного розбиття на кластери. Частка збігів об'єктів в виявлених і вихідних групах є показником ефективності роботи методу;
3. Проблема вибору числа кластерів. Якщо немає ніяких відомостей щодо можливого числа кластерів, необхідно провести ряд експериментів і в результаті перебору різного числа кластерів вибрати оптимальне їх число;
4. Проблема інтерпретації результатів кластеризації. Форма кластерів в більшості випадків визначається вибором методу об'єднання. Проте слід

враховувати, що конкретні методи прагнуть створювати кластери певних форм, навіть якщо в досліджуваному наборі даних кластерів насправді немає [39].

2.3 Обґрунтування вибору методів інтелектуального аналізу для виявлення ознак кіберзагроз у банку

Існує велика кількість методів інтелектуального аналізу даних, які можна застосовувати для побудови моделі виявлення ознак кіберзагроз у банку. Для дослідження даних подій в банківській сфері використаємо: регресії, дерева прийняття рішень, нейронної мережі (табл. 2.2).

Таблиця 2.2 – Методи інтелектуального аналізу

№	Назва методу	Загальна характеристика
1	Логіт-регресія	Різновид множинної регресії, загальне призначення якої полягає в аналізі зв'язку між кількома незалежними змінними і залежною змінною. За допомогою логістичної регресії можна оцінювати ймовірність того, що подія настане для конкретного випробуваного [34].
2	Дерева прийняття рішень	Ідея методу полягає у тому, щоб просуваючись гілками дерева у напрямку справа наліво (тобто від вершини дерева до першої точки прийняття рішення) спочатку розрахувати очікувані виграші по кожній гілці дерева і, порівнюючи ці виграші, зробити остаточний вибір найкращої альтернативи. Мета полягає в тому, щоб створити модель, яка прогнозує значення цільової змінної на основі декількох змінних на вході [19].
3	Нейронні мережі	Математичні моделі, а також їхня програмна та апаратна реалізація, побудовані за принципом функціонування біологічних нейронних мереж – мереж нервових клітин живого організму. Системи, архітектура і принцип дії базується на аналогії з мозком живих істот. Ключовим елементом цих систем виступає штучний нейрон як імітаційна модель нервової клітини мозку людини [32].

Таким чином, методи інтелектуального аналізу, що були обрані для побудови моделей виявлення ознак кіберзагроз в транзакціях користувачів мобільного та інтернет-банкінгу: логіт-регресія, дерево прийняття рішень, нейронна мережа.

2.3.1 Регресія

Завданням дослідження складних систем і процесів часто є перевірка наявності й встановлення типу зв'язку між незалежними змінними x_i (предикторами, факторами), значення яких можуть змінюватися дослідником і мають певну заздалегідь задану похибку, та залежною змінною (відгуком) z . Розв'язання таких завдань є предметом регресійного аналізу. Термін «Регресія» вперше був уведений Ф. Гальтоном наприкінці XIX ст. На практиці завдання регресійного аналізу зазвичай формулюють так: необхідно підібрати достатньо просту функцію, що в певному розумінні найкращим чином описує наявну сукупність емпіричних даних.

Класичний регресійний аналіз включає методи побудови математичних моделей досліджуваних систем, методи визначення параметрів цих моделей і перевірки їх адекватності. Він припускає, що регресія є лінійною комбінацією лінійно незалежних базисних функцій від факторів з невідомими коефіцієнтами (параметрами). Фактори й параметри є детермінованими, а відгуки – рівноточними (тобто мають однакові дисперсії) некорельованими випадковими величинами. Передбачається також, що всі змінні вимірюють у неперервних числових шкалах.

Більш складною проблемою є вибір моделі та її незалежних змінних. У класичному регресійному аналізі припускають, що набір факторів задається однозначно, всі суттєві змінні наявні в моделі й немає ніяких альтернативних способів обрання факторів. На практиці це припущення не виконується. Тому виникає необхідність розробки формальних та неформальних процедур перетворення й порівняння моделей. Для пошуку оптимальних формальних перетворень використовують методи факторного та дискримінантного аналізу. На сьогодні розроблено комп'ютеризовані технології послідовної побудови регресійних моделей.

Важливою особливістю регресійних моделей є те, що їх не можна застосовувати поза межами тієї області значень вихідних параметрів, для якої

вони були побудовані. При використанні регресійних моделей типу полінома, оберненого полінома, тригонометричного ряду та деяких інших слід враховувати, що, збільшуючи кількість членів ряду, можна одержати скільки завгодно близькі до нуля значення функціоналів. Проте це не завжди свідчить про якість апроксимації, оскільки ці функціонали не дають інформації про ступінь наближення моделі до емпіричної залежності у проміжках між наявними точками.

Поліноміальні регресійні моделі, як правило, є формальними. Їх використовують для опису систем і процесів, теорію яких розроблено недостатньо. При цьому спираються на відомі властивості ряду Тейлора для аналітичних функцій. Більш цікавими для дослідників зазвичай є математичні моделі, які відображають структуру та зв'язки у системах, сутність і механізми процесів, що відбуваються у них. Якщо теоретичні основи досліджуваних систем і процесів достатньо розроблені, часто постає проблема визначення окремих параметрів моделі за наявними емпіричними даними. Для її вирішення у багатьох випадках можна використовувати формальні процедури регресійного аналізу [15].

Головна особливість регресійного аналізу: за допомогою нього можна отримати конкретні відомості про те, яку форму та характер має залежність між досліджуваними змінними.

Послідовність етапів регресійного аналізу:

1. Формулювання задачі. На цьому етапі формуються попередні гіпотези про залежність досліджуваних явищ.
2. Визначення залежних та незалежних (пояснюючих) змінних.
3. Збір статистичних даних. Дані повинні бути зібрані для кожної зі змінних, включених в модель.
4. Формулювання гіпотези про форму зв'язку (простий чи множинний, лінійний чи нелінійний).

5. Визначення функції регресії (полягає в розрахунку чисельних значень параметрів рівняння регресії).

6. Оцінка точності регресійного аналізу.

7. Інтерпретація отриманих результатів. Отримані результати регресійного аналізу порівнюються з попередніми гіпотезами. Оцінюється коректність та правдоподібність отриманих результатів.

8. Передбачення невідомих значень залежної змінної.

За допомогою регресійного аналізу можливе вирішення завдання прогнозування і класифікації. Прогнозні значення обчислюються шляхом підстановки в рівняння регресії параметрів значень пояснюючих змінних. Рішення завдання класифікації здійснюється таким чином: лінія регресії ділить всю множину об'єктів на два класи, і та частина множини, де значення функції більше нуля, належить до одного класу, а та частина, де воно менше нуля, – до іншого класу.

Основні завдання регресійного аналізу: встановлення форми залежності, визначення функції регресії, оцінка невідомих значень залежної змінної.

Встановлення форми залежності. Характер і форма залежності між змінними можуть утворювати такі різновиди регресії: позитивна лінійна регресія (виражається в рівномірному зростанні функції); позитивна рівноприскорено зростаюча регресія; позитивна рівноуповільнено зростаюча регресія; негативна лінійна регресія (виражається в рівномірному падінні функції); негативна рівноприскорено спадна регресія; негативна рівноуповільнено спадна регресія.

Визначення функції регресії. Це завдання зводиться до з'ясування дії на залежну змінну головних чинників або причин при незмінних інших рівних умовах і за умови виключення впливу на залежну змінну випадкових елементів. Функція регресії визначається у вигляді математичного рівняння того або іншого типу.

Оцінка невідомих значень залежної змінної. Вирішення цього завдання зводиться до вирішення задачі одного з типів:

- оцінка значень залежної змінної всередині розглянутого інтервалу вихідних даних, тобто пропущених значень; при цьому вирішується завдання інтерполяції;
- оцінка майбутніх значень залежної змінної, тобто знаходження значень поза заданого інтервалу вихідних даних; при цьому вирішується завдання екстраполяції.

Обидва завдання вирішуються шляхом підстановки в рівняння регресії знайдених оцінок параметрів значень незалежних змінних. Результат рішення рівняння являє собою оцінку значення цільової (залежної) змінної [39].

В логістичній регресійній моделі змодельовані значення залежної змінної знаходяться в інтервалі від 0 до 1 незалежно від значень незалежних змінних, тому, ця модель часто використовується для аналізу бінарних залежних змінних або змінних відгуку.

При цьому ймовірність настання події визначається функцією (2.1):

$$\text{logit}(\hat{p}) = \log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{w}_0 + \hat{w}_1 \cdot x_1 + \hat{w}_2 \cdot x_2 + \dots + \hat{w}_n \cdot x_n, \quad (2.1)$$

де \hat{p} – змодельована ймовірність настання події, що цікавить дослідника;

x_n – n -ий незалежний фактор;

w_n – n -ий коефіцієнт регресії (оцінка фактору);

n – порядковий номер фактору.

Однак, застосування логістичного перетворення до рівняння логіт-регресії породжує певні проблеми. При вирішенні задачі звичайної лінійної регресії до спостережуваних значень підганяється деяка гіперповерхня – пряма у випадку простої регресії, площина – у випадку двох незалежних змінних. Також вимогою є нормальність і некорельованість помилок.

Тому для оцінки параметрів логіт-регресії використовується тільки метод максимальної правдоподібності, за якого процес оцінки регресійних коефіцієнтів зводиться до максимізації ймовірності появи конкретної вибірки (при заданих спостережуваних значеннях). Це часто призводить до невисокого відсотку коректної класифікації. Логіт-регресія також слабо стійка до надмірної підгонки даних [27].

Для вирішення даних завдань моделювання також можна використовувати пробіт модель – це статистична модель бінарного вибору, що використовується для передбачення ймовірності виникнення події на основі функції стандартного нормального розподілу. Модель пробіт регресії, також як і модель логістичної регресії, відносять до моделей бінарного вибору, тому функції і завдання її побудови аналогічні логіт моделі.

У моделі пробіт регресії розрахункове значення залежної змінної виражається як значення функції розподілу стандартного нормального закону. Пробіт – це значення, для якого обчислюється функція розподілу стандартного нормального закону розподілу. Значення пробіта залежить від лінійних комбінацій значень факторних змінних. Як і для логіт моделі, залежна змінна в пробіт моделі є бінарною. Фактори в пробіт моделі можуть бути кількісними змінними або категоріальними, перетвореними в бінарні змінні.

Модель бінарного вибору називається пробіт регресією, якщо вона задовольняються такі дві умови:

- 1) залишки моделі бінарного вибору є випадковими нормально розподіленими величинами;
- 2) функція розподілу ймовірностей є нормальною ймовірнісною функцією.

Пробіт регресія може бути представлена формулами (2.2 – 2.3):

$$p_i \left(y_i = \frac{1}{x_1 \cdots x_n} \right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-s^2/2} ds; \quad (2.2)$$

$$z = \hat{w}_0 + \hat{w}_1 \cdot x_1 + \hat{w}_2 \cdot x_2 + \cdots + \hat{w}_n \cdot x_n, \quad (2.3)$$

де p_i – змодельована ймовірність настання події, що цікавить дослідника;

x_n – n -ий незалежний фактор;

\hat{w}_n – n -ий коефіцієнт регресії (оцінка фактору);

n – порядковий номер фактору.

Для оцінки параметрів, як і у логіт моделі, використовується метод максимальної правдоподібності. Значної різниці в результатах розрахунку пробіт чи логіт моделі немає. Логістична функція щільності дуже близька до нормального стандартного розподілу, але з більш «товстими хвостами». Різниця у використанні цих моделей можлива лише у випадку, якщо дані концентруються у «хвостах» розподілів. Але історично в моделюванні соціальних явищ використовується частіше логіт, адже за допомогою нього можна розрахувати більшу кількість показників, наприклад: шанси настання певної події, їх відношення [20].

2.3.2 Дерево рішень

Дерева рішень (Decision Trees) – це метод, що дозволяє передбачати приналежність спостережень або об'єктів до того чи іншого класу категоріальної залежної змінної в залежності від відповідних значень однієї або декількох предикторних змінних. Мета побудови дерев рішень полягає в передбаченні (або поясненні) значень категоріальної залежної змінної, і тому використовувані методи тісно пов'язані з більш традиційними методами дискримінантного аналізу, кластерного аналізу, непараметричної статистики.

Широка сфера застосування дерев рішень робить їх вельми привабливим інструментом аналізу даних, але не слід тому вважати, що його

рекомендується використовувати замість традиційних методів статистики. Навпаки, якщо виконані більш строгі теоретичні припущення, що накладаються традиційними методами, і вибірковий розподіл має деякі спеціальні властивості, то більш результативним буде використання саме традиційних методів [39].

В найбільш простому вигляді дерево рішень – це спосіб представлення правил в ієрархічній, послідовній структурі. Основа такої структури – відповіді «Так» чи «Ні» на ряд питань.

Зазвичай кожен вузол дерева включає перевірку певної незалежної змінної. Іноді в вузлі дерева дві незалежні змінні порівнюються одна з одною або визначається деяка функція від однієї або декількох змінних.

Якщо змінна, яка перевіряється в вузлі, приймає категоріальні значення, то кожному можливому значенню відповідає гілка, що виходить з вузла дерева. Якщо значенням змінної є число, то перевіряється, більше або менше це значення деякої константи. Іноді область числових значень розбивають на декілька інтервалів. В цьому випадку виконується перевірка на потрапляння значення в один з інтервалів.

Листя дерев відповідають значенням залежної змінної, тобто класам. Об'єкт належить певному класу, якщо значення його незалежних змінних задовольняють умовам, записаним в вузлах дерева на шляху від кореня до листа, відповідному цього класу.

Якщо будь-яка незалежна змінна об'єкта, що класифікується не має значення, то виникає проблема, пов'язана з невизначеністю шляху, за яким необхідно рухатися по дереву. У деяких випадках пропущені значення можна замінювати значеннями за замовчуванням. Якщо такий підхід неприйнятний, то необхідно передбачити спеціальні способи обробки таких ситуацій (наприклад, переміщатися по гілці, яка веде до більшої кількості об'єктів з навчальної вибірки). Інший варіант обробки може бути пов'язаний з додаванням спеціальної гілки до вузла для пропущених значень.

Дерева рішень легко перетворюються в правила. В умовну частину таких правил записується умова, описане в вузлах дерева на шляху до листу, в заключну частину-значення, визначене в листі [14].

Правило або способи розбивки множин записів або варіантів називають вирішальним правилом (2.4):

$$a_{i,k} = \begin{cases} 1 \\ 0 \end{cases}, \quad (2.4)$$

де $a_{i,k} = 1$, якщо умова s_i для правила r_k виконується; 0, в іншому випадку;

$S\{s_i\}, i = \overline{1, l}$ – множина умов, що описують параметри обраної предметної області;

$R\{r_k\}, k = \overline{1, m}$ – множина вирішальних правил, що описують конкретні дії, що виконуються при заданих значеннях параметрів з множини умов.

Це правило фактично є логічною структурою «якщо ..., то ...», що ділить аналізовану множину на дві групи. По мірі опускання по дереву рішень від вершини до листків, створюється усе більше відфільтрованих однорідних множин, що задовольняють певному набору умов, сформульованих у вузлах дерева.

Дерева рішень мають ряд очевидних переваг перед статистичними методами.

Інтуїтивність дерев рішень. Класифікаційна модель, представлена у вигляді дерева рішень, є інтуїтивною і спрощує розуміння розв'язуваної задачі. Результат роботи алгоритмів конструювання дерев рішень, на відміну, наприклад, від нейронних мереж, що представляють собою «чорні ящики», легко інтерпретується користувачем. Ця властивість дерев рішень не тільки важлива при віднесенні до певного класу нового об'єкта, а й корисна при інтерпретації моделі класифікації в цілому. Дерево рішень дозволяє

зрозуміти і пояснити, чому конкретний об'єкт відноситься до того чи іншого класу. Дерева рішень дають можливість отримувати правила з бази даних на природній мові.

Рішення важко формалізованих завдань. Дерева рішень дозволяють створювати класифікаційні моделі в тих областях, де аналітику досить складно формалізувати знання. Алгоритм конструювання дерева рішень не вимагає від користувача вибору вхідних атрибутів (незалежних змінних). На вхід алгоритму можна подавати всі існуючі атрибути, алгоритм сам вибере найбільш значущі серед них, і тільки вони будуть використані для побудови дерева. У порівнянні, наприклад, з нейронними мережами, це значно полегшує користувачеві роботу, оскільки в нейронних мережах вибір кількості вхідних атрибутів істотно впливає на час навчання.

Точність моделі. Точність моделей, створених за допомогою дерев рішень, порівнянна з іншими методами побудови класифікаційних моделей (статистичні методи, нейронні мережі).

Швидкий процес навчання. На побудову класифікаційних моделей за допомогою алгоритмів конструювання дерев рішень потрібно значно менше часу, ніж, наприклад, на навчання нейронних мереж.

Використання категоріальних видів даних. Більшість алгоритмів конструювання дерев рішень мають можливість спеціальної обробки пропущених значень. Багато класичних статистичних методів, за допомогою яких вирішуються завдання класифікації, можуть працювати тільки з числовими даними, в той час як дерева рішень працюють і з числовими, і з категоріальними типами даних. Багато статистичних методів є параметричними, і користувач повинен заздалегідь володіти певною інформацією, наприклад, знати вид моделі, мати гіпотезу про вид залежності між змінними, припускати, який вид розподілу мають дані. Дерева рішень, на відміну від таких методів, будують непараметричні моделі. Таким чином,

дерева рішень здатні вирішувати такі завдання Data Mining, в яких відсутня апріорна інформація про вид залежності між досліджуваними даними [39].

2.3.3 Нейронна мережа

Робота над штучними нейронними мережами, які зазвичай називають «нейронними мережами», була мотивована від самого початку, визнаючи, що людський мозок обчислює цілком іншим шляхом, ніж звичайний цифровий комп'ютер. Мозок – це дуже складний, нелінійний і паралельний комп'ютер (система обробки інформації). Він має можливість організувати структурні компоненти, відомі як нейрони, для того, щоб виконувати певні обчислення (наприклад, розпізнавання образів, сприйняття та управління двигуном) у багато разів швидше, ніж найшвидший цифровий комп'ютер, який існує сьогодні.

Нейронна мережа – це масовий паралельно розподілений процесор, що складається з простих блоків обробки, які мають природну схильність зберігати експериментальні знання та зробити їх доступним для використання. Він схожий на мозок у двох аспектах:

- 1) знання набуває мережа від навколишнього середовища через процес навчання;
- 2) сильні сторони зв'язку нейронів, відомі як синаптичні ваги, використовуються для зберігання отриманих знань.

Процедура, яка використовується для виконання навчального процесу, називається алгоритмом навчання, функцією якого є правильна модифікація синаптичних вагів мережі для досягнення бажаної проектної мети [11].

Нейронні мережі – це клас моделей, заснованих на біологічній аналогії з мозком людини і призначених (після проходження етапу так званого навчання на наявних даних) для вирішення різноманітних завдань аналізу даних. При застосуванні цих методів перш за все постає питання вибору конкретної архітектури мережі (числа «шарів» і кількості «нейронів» в

кожному з них). Розмір і структура мережі повинні відповідати (наприклад, в сенсі формальної обчислювальної складності) суті досліджуваного явища. Оскільки на початковому етапі аналізу природа явища зазвичай відома погано, вибір архітектури є непростим завданням і часто пов'язаний з тривалим процесом «проб і помилок» (проте останнім часом стали з'являтися нейронно-мережеві програми, в яких для вирішення трудомісткого завдання пошуку кращої архітектури мережі застосовуються методи штучного інтелекту) [33].

Математична модель для вирішення завдань машинного навчання реалізується групою з'єднаних нейронів для моделювання нелінійних залежностей.

Функція активації в нейронній мережі – це функція, що обчислює розраховує значення вихідного шару нейрона Існує багато різних функцій активації, наприклад: функція одиничного скачку, кусково-лінійна, сигмоїдальна тощо. Їх використовуються в залежності від того, на якому числовому інтервалу має належати прогнозована величина.

Бінарну прогнозну формулу нейронної мережі та формулу n -го нейрона наведено у формулах (2.5) та (2.6), відповідно [40].

$$\hat{y} = \hat{w}_{00} + \hat{w}_{01} \cdot H_1 + \hat{w}_{02} \cdot H_2 + \dots + \hat{w}_{0n} \cdot H_n; \quad (2.5)$$

$$H_n = \tanh(\hat{w}_{n0} + \hat{w}_{n1} \cdot x_1 + \hat{w}_{n2} \cdot x_2 + \dots + \hat{w}_{nk} \cdot x_k) \quad (2.6)$$

де \hat{p} – змодельована ймовірність настання події, що цікавить дослідника;

w_n – n -та вагова оцінка;

H_n – n -ий прихований елемент;

n – порядковий номер елементу;

x_k – k -ий незалежний фактор;

k – порядковий номер фактору.

Розрізняють два типи нейронних мереж:

- мережі прямого поширення (Feed forward Neural Networks);
- рекурентні нейронні мережі (Recurrent Neural Networks).

В мережах прямого поширення сигнал передається від вхідного рівня нейронів до вихідного по «прошаркам». Відбувається розрахунок нелінійних вихідних функцій, від вхідних змінних кожна, як композиції алгебраїчних функцій активації. Немає затримок, часу, тому що немає циклів.

В рекурентних нейронних мережах присутні довільні топології з циклами. Відбувається моделювання системи з станами (динамічні системи). Є поняття «затримки» у деяких вагів. Процес навчання – важкий, а результат не завжди передбачуваний: нестабільний (нестійкий) сигнал на виході, несподівана поведінка (осциляції, хаос) [37].

Використання нейронних мереж забезпечує наступні корисні властивості систем.

Нелінійність. Штучні нейрони можуть бути лінійними і нелінійними. Нейронні мережі, побудовані із з'єднань нелінійних нейронів, самі є нелінійними. Більше того, ця нелінійність особливого сорту, оскільки вона розподілена по мережі. Нелінійність є надзвичайно важливою властивістю, особливо якщо сам фізичний механізм, відповідальний за формування вхідного сигналу, теж є нелінійним (наприклад, людська мова).

Відображення вхідної інформації у вихідну. Однією з популярних парадигм навчання є навчання з учителем. Це має на увазі зміну синаптичних ваг на основі набору маркованих навчальних прикладів. Кожен приклад складається з вхідного сигналу і відповідного йому бажаного відгуку. З цієї безлічі випадковим чином вибирається приклад, а нейронна мережа модифікує синаптичні ваги для мінімізації розбіжностей бажаного вихідного сигналу і формованого мережею відповідно до обраного статистичному критерію. При цьому власне модифікуються вільні параметри мережі. Раніше використані приклади можуть згодом бути застосовані знову, але вже в

іншому порядку. Це навчання проводиться до тих пір, поки зміни синаптичних ваг не стануть незначними. Таким чином, нейронна мережа навчається на прикладах, складаючи таблицю відповідностей вхід-вихід для конкретного завдання.

Адаптивність. Нейронні мережі мають здатність адаптувати свої синаптичні ваги до змін навколишнього середовища. Зокрема, нейронні мережі, навчені діяти в певному середовищі, можуть бути легко перевчити для роботи в умовах незначних коливань параметрів середовища. Більш того, для роботи в нестаціонарному середовищі можуть бути створені нейронні мережі, що змінюють синаптичні ваги в реальному часі. Природна для класифікації образів, обробки сигналів і завдань управління архітектура нейронних мереж може бути об'єднана з їх здатністю до адаптації, що призведе до створення моделей адаптивної класифікації образів, адаптивної обробки сигналів і адаптивного керування. Відомо, що чим вище адаптивні здібності системи, тим більш стійкою буде її робота в нестаціонарному середовищі.

Контекстна інформація. Знання представляються в самій структурі нейронної мережі за допомогою її стану активації. Кожен нейрон мережі потенційно може бути підданий впливу всіх інших її нейронів. Як наслідок, існування нейронної мережі безпосередньо пов'язане з контекстною інформацією.

Аналогія з нейробіологією. Будова нейронних мереж визначається аналогією з людським мозком, який є живим доказом того, що відмовостійкі паралельні обчислення не тільки фізично реалізовані, але і є швидким і потужним інструментом вирішення завдань. Нейробіологи розглядають штучні нейронні мережі як засіб моделювання фізичних явищ. З іншого боку, інженери постійно намагаються почерпнути у нейробіологів нові ідеї, що виходять за рамки традиційних електросхем [32].

З точки зору машинного навчання, нейронна мережа являє собою окремий випадок методів розпізнавання образів, дискримінантного аналізу, методів кластеризації тощо. З математичної точки зору, навчання нейронних мереж – це багатопараметрична задача нелінійної оптимізації. З точки зору кібернетики, нейронна мережа використовується в задачах адаптивного управління і як алгоритми для робототехніки. З точки зору розвитку обчислювальної техніки та програмування, нейронна мережа – спосіб вирішення проблеми ефективного паралелізму.

Нейронні мережі не програмуються в звичайному розумінні цього слова, вони навчаються. Можливість навчання – одна з головних переваг нейронних мереж перед традиційними алгоритмами. Технічно навчання полягає в знаходженні коефіцієнтів зв'язків між нейронами. У процесі навчання нейронна мережа здатна виявляти складні залежності між вхідними даними і вихідними, а також виконувати узагальнення. Це означає, що у разі успішного навчання мережа зможе повернути вірний результат на підставі даних, які були відсутні в навчальній вибірці, а також неповних та/або «зашумлених», частково перекручених даних.

Після того, як нейронна мережа навчена, можна застосовувати її для вирішення необхідних завдань. Нейронна мережа, коректно навчена, може з великою ймовірністю правильно реагувати на нові, невідомі їй раніше дані.

РОЗДІЛ 3 ПРАКТИЧНА РЕАЛІЗАЦІЯ ІНТЕЛЕКТУАЛЬНОГО АНАЛІЗУ КІБЕРЗАГРОЗ ІЗ ВИКОРИСТАННЯМ АНАЛІТИЧНОГО ПАКЕТУ SAS ENTERPRICE MINER

3.1 Опис програмної реалізації модельних розрахунків на EOM

Для побудови моделі виявлення кібернетичних загроз в банківських установах з метою майбутнього попередження цих загроз в разі їх виникнення скористаємося аналітичним пакетом SAS Enterprise Miner.

SAS Enterprise Miner полегшує і систематизує процес інтелектуального аналізу даних, дозволяючи створювати високоточні передбачувальні і описові моделі на основі аналізу величезної кількості інформації, що збирається у всій організації. Цей пакет інструментів допомагає вирішувати широке коло завдань, що вимагають вивчення інформації і можливості передбачити хід подій, а саме: виявляти випадки шахрайства, визначати і мінімізувати рівень ризиків, прогнозувати потреби в ресурсах, попереджати інциденти, підвищувати рівень відгуку на маркетингові кампанії, знижувати відтік клієнтів та інші.

Цей пакет являє собою найбільш потужне і повнофункціональне рішення з усіх наявних на ринку для передбачувальної аналітики та інтелектуального аналізу даних. SAS Enterprise Miner дозволяє користувачам досліджувати і аналізувати складні дані, знаходити стійкі закономірності і, ґрунтуючись на фактах і отриманих висновках, приймати виважені рішення.

SAS Enterprise Miner створений для фахівців з аналізу даних, статистиків, маркетингових аналітиків, маркетологів, експертів з аналізу ризиків, фахівців з виявлення шахрайських дій. Цей інструмент також активно використовується інженерами, науковцями та бізнес-аналітиками, яким необхідно розуміти і аналізувати постійно зростаючі обсяги даних,

розпізнавати критичні завдання бізнесу або наукових досліджень і приймати обґрунтовані рішення [10].

3.1.1 Проведення первинного та кластерного аналізу змінних

Для реалізації поставленої задачі відкриємо програму SAS Enterprise Miner та створимо новий проект з ім'ям diploma, виконавши наступні кроки.

New Project > Next. Введемо ім'я проекту diploma. Вкажемо папку для збереження проекту. Next > Next > Finish.

В створеному проекті виконаємо підключення бібліотеки Diploma та створимо діаграму з ім'ям Bank.

File > New diagram > Name = Bank.

Задамо джерело даних banking.sas7bdat.

File > New > Data Source > Next > Browse > banking.sas7bdat > Next.

Додана вибірка даних містить 200000 записів і 8 параметрів.

На кроці 4 Metadata Advisor Options натиснемо Advanced > Customize > змінимо значення властивостей та натиснемо Next (рис. 3.1).

Class Levels Count Threshold = 2, означає, що тільки бінарні чисельні змінні будуть сприйматися як категоріальні. А всі інші чисельні змінні у яких більш ніж два рівня будуть сприйняті як інтервальні (безперервні).

Reject Levels Count Threshold = 100, означає, що змінні не будуть відхилені з аналізу через велику кількість рівнів.

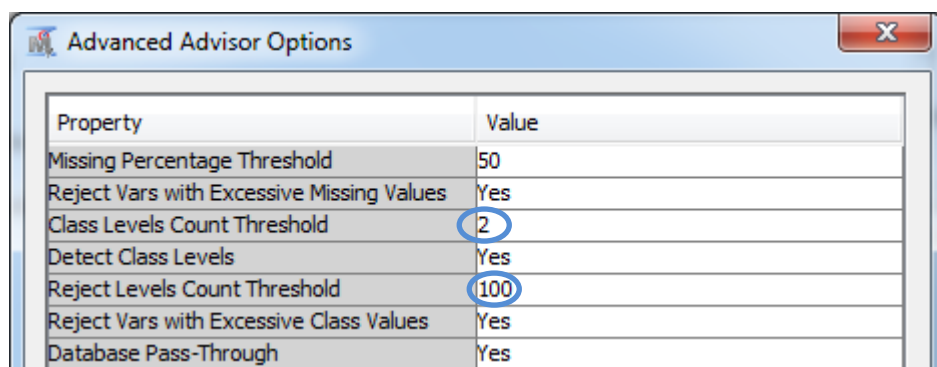


Рисунок 3.1 – Додаткові налаштування

Для цільової змінної з набору даних isfraud, яка відповідає за відгук, задамо роль Target, рівень – Binary (рис. 3.2). Завдяки цьому система автоматично обере логістичну регресію):

1 – так (виконана транзакція є загрозою);

0 – ні.

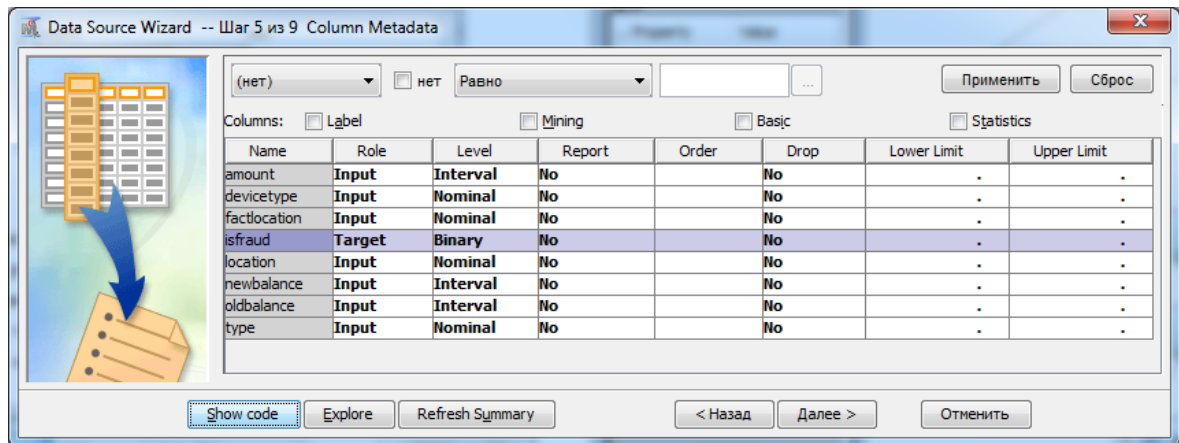


Рисунок 3.2 – Визначення ролей вхідних змінних

Для завершення створення джерела даних обираємо Next > Next > Next > Finish.

Перш ніж проводити інтелектуальний аналіз, виконаємо первинний аналіз вхідних даних за допомогою інструмента StatExplore пакету SAS Enterprise Miner.

Перетягнемо джерело даних BANKING у вікно робочої області Bank. Додамо інструмент StatExplore та об'єднаємо з джерелом даних (рис. 3.3).

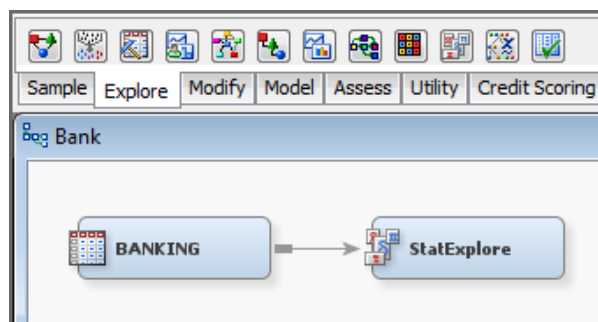


Рисунок 3.3 – Об'єднання інструмента StatExplore з джерелом даних

Натиснемо правою кнопкою по вузлу StatExplore і виберемо Run з меню швидкого виклику. Переглянемо результати ходу виконання даного вузла.

На рисунку 3.4 відображені категоріальні змінні та їх основні властивості: роль змінної, кількість рівнів, пропущенні значення, мода.

Output									
36	Data Role=TRAIN								
37									
38									
39	Data	Variable		Number			Mode		Mode2
40	Role	Name	Role	Levels	Missing	Mode	Percentage	Mode2	Percentage
41									
42	TRAIN	devicetype	INPUT	2	0	M	50.53	I	49.47
43	TRAIN	factlocation	INPUT	2	0	UA	79.12	Other	20.88
44	TRAIN	type	INPUT	5	0	PAYMENT	39.51	CASH_OUT	30.72
45	TRAIN	isfraud	TARGET	2	0	0	79.12	1	20.88

Рисунок 3.4 – Основні властивості вхідних категоріальних змінних

На рисунку 3.5 відображена інформація стосовно цільової змінної isfraud: частоти позитивного та негативного відгуку, а також долі від цілого.

Доля проведених банківських транзакцій, які виявились кібернетичними загрозами становить 20,9 %, в свою чергу в 79,1% проведених операцій не виявлено кіберзагроз.

Output						
52	Data Role=TRAIN					
53						
54	Data	Variable			Frequency	
55	Role	Name	Role	Level	Count	Percent
56						
57	TRAIN	isfraud	TARGET	0	79124	79.124
58	TRAIN	isfraud	TARGET	1	20876	20.876

Рисунок 3.5 – Статистична інформація щодо цільової змінної isfraud

На рисунку 3.6 відображена статистична інформація по інтервальних змінних: роль змінної, середнє значення, стандартне відхилення, пропущені значення, мінімум, медіана, максимум.

Output												
65	Data Role=TRAIN											
66												
67												
68	Variable	Role	Mean	Standard	Non	Missing	Missing	Minimum	Median	Maximum	Skewness	Kurtosis
69				Deviation	Missing							
70	amount	INPUT	187512.4	478553.1	99962	38	0	54609	33966807	19.2124	877.3581	
71	newbalance	INPUT	661149.4	2386766	98091	1909	0	0	99696007	16.28264	486.3307	
72	oldbalance	INPUT	652039.6	2365850	98165	1835	0	18545	99696007	16.56692	501.0274	

Рисунок 3.6 – Статистичні характеристики вхідних інтервальних змінних

В результаті проведеного первинного аналізу було отримано основні статистичні характеристики вхідних змінних, визначено ролі змінних у моделюванні, а також виявлено, що у вхідному масиві даних відсутні пропущені значення в інтервальних змінних.

Для розбиття набору даних на тренувальний, тестовий та перевірочний набори даних скористаймося інструментом Data Partition пакету SAS Enterprise Miner.

Додамо даний інструмент та об'єднаємо з джерелом даних (рис. 3.7).

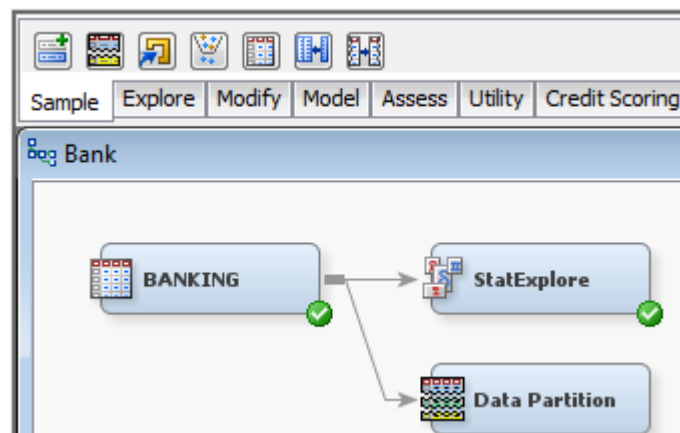


Рисунок 3.7 – Додавання інструмента Data Partition в робочу область

У властивостях вузла Data Partition оберемо частки даних для навчання (50%) та перевірки (50%) як вказано на рисунку 3.8.

.. Property	Value
Train	
Variables	
Output Type	Data
Partitioning Method	Default
Random Seed	12345
Data Set Allocations	
Training	50.0
Validation	50.0
Test	0.0

Рисунок 3.8 – Налаштування властивостей вузла Data Partition

Далі, проаналізувавши графіки інтервальних змінних, можна побачити, що розподіл даних величин не відповідає нормальному закону розподілу (рис. 3.9).

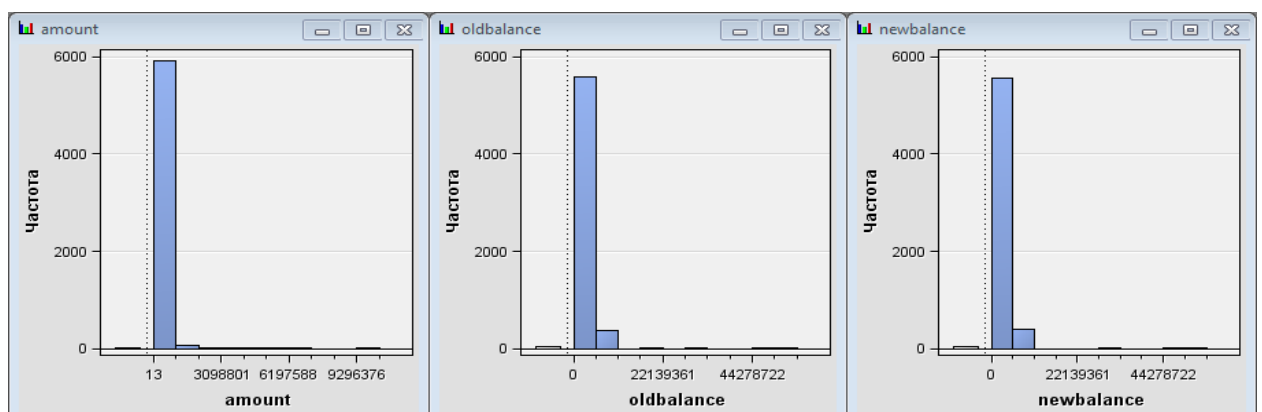


Рисунок 3.9 – Перевірка нормального закону розподілу у вхідних інтервальних змінних

А тому, для подальшої побудови моделей необхідно прологіритмувати вхідні змінні. Для цього скористаймося інструментом Transform Variables пакету SAS Enterprise Miner.

Додамо у вікно робочої області інструмент Transform Variables та об'єднаємо з вузлом Data Partition (рис. 3.10).

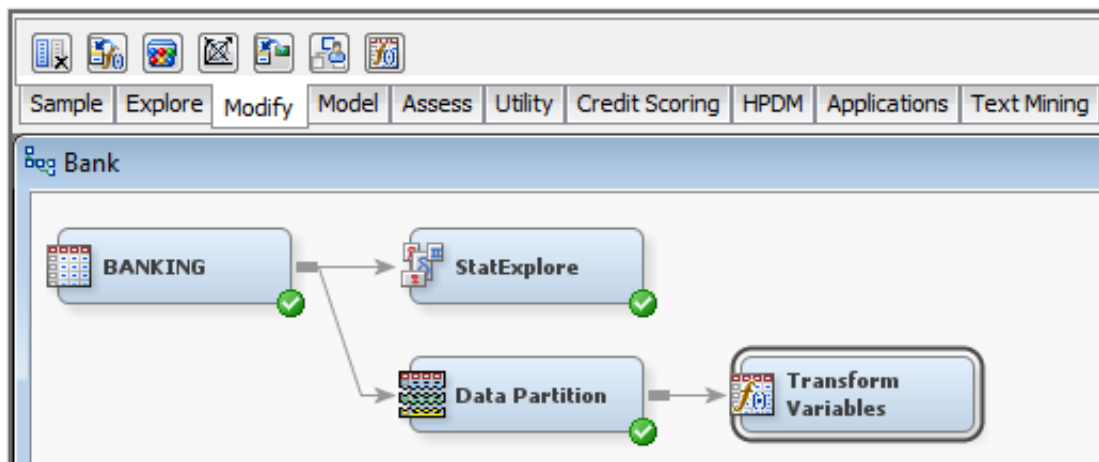


Рисунок 3.10 – Додавання вузла Transform Variables в робочу область

У властивостях вузла Transform Variables оберемо Variables та вкажемо для інтервальних змінних amount, oldbalance та newbalance метод Log (рис. 3.11).

Variables - Trans2

(нет) ☐ нет Равно

Columns: ☐ Label ☐ Mining

Name	Method	Number of Bins	Role	Level
amount	Log	4	Input	Interval
devicetype	Default	4	Input	Nominal
factlocation	Default	4	Input	Nominal
isfraud	Default	4	Target	Binary
location	Default	4	Input	Nominal
newbalance	Log	4	Input	Interval
oldbalance	Log	4	Input	Interval
type	Default	4	Input	Nominal

Рисунок 3.11 – Логарифмування вхідних інтервальних змінних

Після проведення первинного аналізу даних та логарифмування вхідних змінних, джерело даних можна застосовувати для інтелектуального аналізу даних.

Для виявлення прихованих, неочевидних тенденцій та закономірностей у вхідних даних, проведемо більш серйозний, глибинний статистичний аналіз – кластерний. Дослідження виконаємо у пакеті SAS Enterprise Miner.

Спочатку обираємо вхідні змінні для кластерного аналізу. Вхідні змінні повинні мати наступні властивості:

- бути значимими для цілей аналізу;
- бути відносно незалежними;
- бути обмеженими по кількості [25].

Зважаючи на ці вимоги, було обрано наступні вхідні змінні (табл. 3.1).

Таблиця 3.1 – Опис вхідних змінних

Ім'я змінної	Економічний зміст	Роль змінної	Тип
amount (X_1)	Загальна сума, що була проходила в транзакціях	вхідна	interval
devicetype (X_2)	Тип пристрою, з якого виконувалась транзакція	вхідна	nominal
factlocation (X_3)	Зафіксоване місцеположення пристрою, з якого проводилась транзакція	вхідна	nominal
newbalance (X_5)	Баланс клієнта після проведення транзакції	вхідна	interval
type (X_7)	Тип виконаної транзакції	вхідна	nominal

Додамо в область діаграми інструмент Cluster та об'єднаємо з вузлом Transform Variables (рис. 3.12).

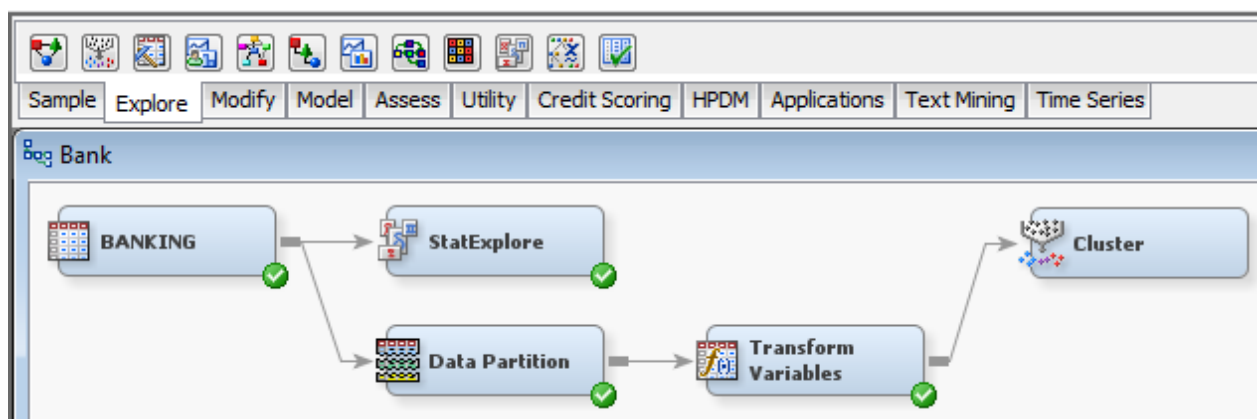


Рисунок 3.12 – Додавання інструмента Cluster в робочу область діаграми

У властивостях вузла Cluster вкажемо самотійно максимальну кількість кластерів – 4 (рис. 3.13).

Property	Value
Train	
Variables	
Internal Standardization	Standardization
<input checked="" type="checkbox"/> Number of Clusters	
Specification Method	User Specify
Maximum Number of Clusters	4
<input checked="" type="checkbox"/> Selection Criterion	
Clustering Method	Ward
Preliminary Maximum	50
Minimum	2
Final Maximum	20
CCC Cutoff	3

Рисунок 3.13 – Налаштування властивостей вузла Cluster

Результатом кластерного аналізу даних у пакеті SAS Enterprise Miner є виділення 4-х кластерів з наступними статистичними характеристиками (табл. 3.2).

Таблиця 3.2 – Статистика у розрізі окремих кластерів в пакеті SAS Enterprise Miner

Характеристики	№ сегмента кластеру			
	1	2	3	4
Кількість випадків, що потрапили у кластер	48026	43793	60528	47653
Відсоток випадків, що потрапили у кластер	24,01	21,9	30,26	23,83
Найближчий кластер до даного	4	3	1	1
Середнє значення LOG_amount у кластері	8,5590	11,7833	11,1866	11,8361
Середнє значення LOG_newbalance у кластері	10,0476	0,0002	0,0087	13,5509

Діаграму розподілу даних по кластерам представлено на рис. 3.14.

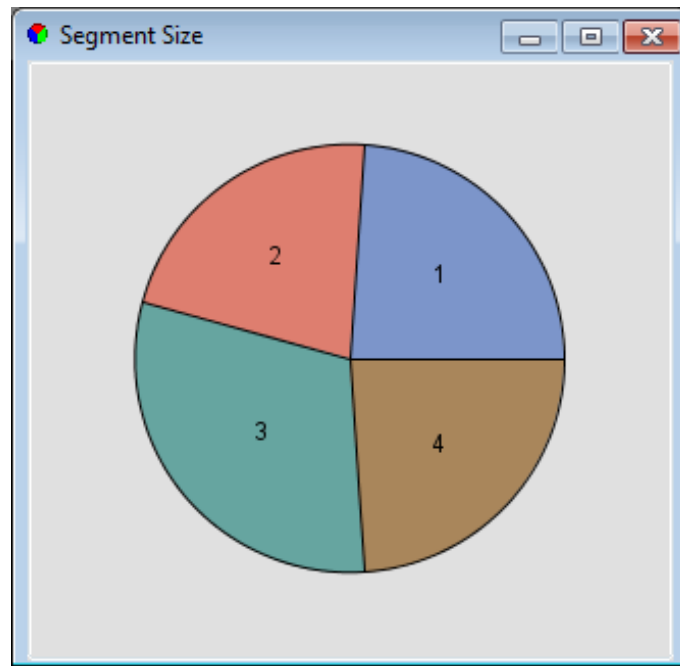


Рисунок 3.14 – Розподіл даних на кластери в пакеті SAS Enterprise Miner

Отже, кількість випадків, що класифіковано у 1-й кластер – 48026, у 2-й – 43793, у 3-й – 60528, у 4-й – 47653. Тобто, за величиною випадків кластери є приблизно однаковими.

Оскільки для генерування сегментів використовується більше трьох змінних, інтерпретація таких графіків стає складнішою. Для цього в SAS Enterprise Miner є інструмент для інтерпретації композиції кластерів: Segment Profile на панелі Assess, який дозволяє порівнювати розподіл змінної в конкретному сегменті з розподілом змінної в загальному наборі даних. Також, змінні упорядковуються відносно того, наскільки добре вони характеризують даний сегмент. Додаємо інструмент Segment Profile з набору інструментів Assess в робочу область діаграми та з'єднуємо його з вузлом Cluster для дослідження кожного кластеру окремо (рис.3.15).

Запускаємо на виконання вузол вузол Segment Profile і обираємо Results. Відкриється вікно Results. Профільне дослідження кластерів та важливість кожної змінної у формуванні того чи іншого кластеру наведені на рисунках 3.16 – 3.17.

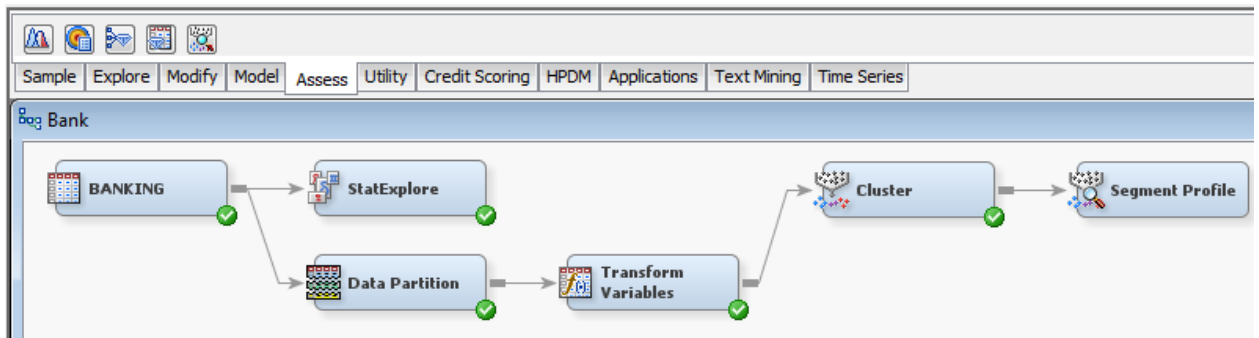


Рисунок 3.15 – Додавання інструменту Segment Profile в область діаграми

Таким чином, при формуванні першого кластеру найбільшу вагу мали змінні LOG_amount та LOG_newbalance, незначний вплив становила змінна factlocation. На формування другого кластеру найбільше вплинула змінна factlocation та менш значно вплинули змінні LOG_newbalance та LOG_amount. Змінна LOG_newbalance спричинила значний вплив на формування третього та четвертого кластерів, в той час як вплив змінних LOG_amount та factlocation на ці кластери був меншим.

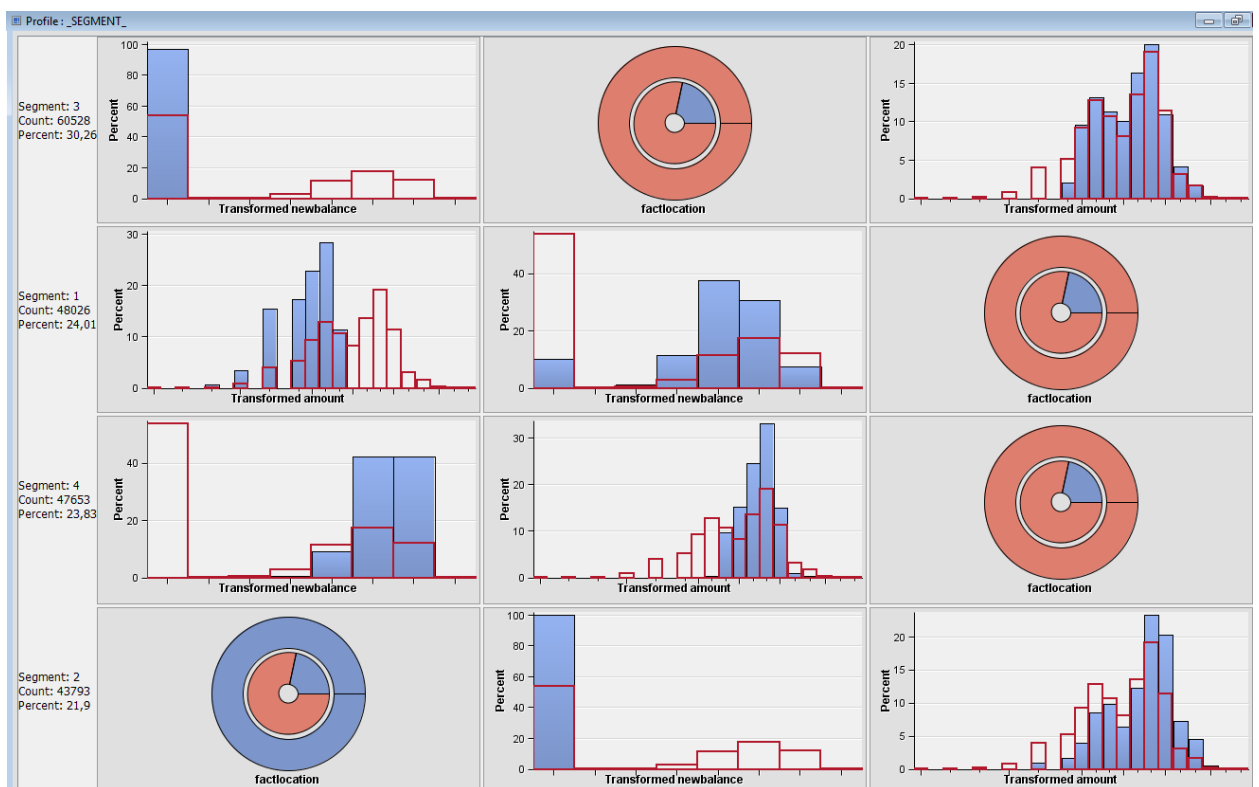


Рисунок 3.16 – Профільний аналіз кластерів в пакеті SAS Enterprise Miner

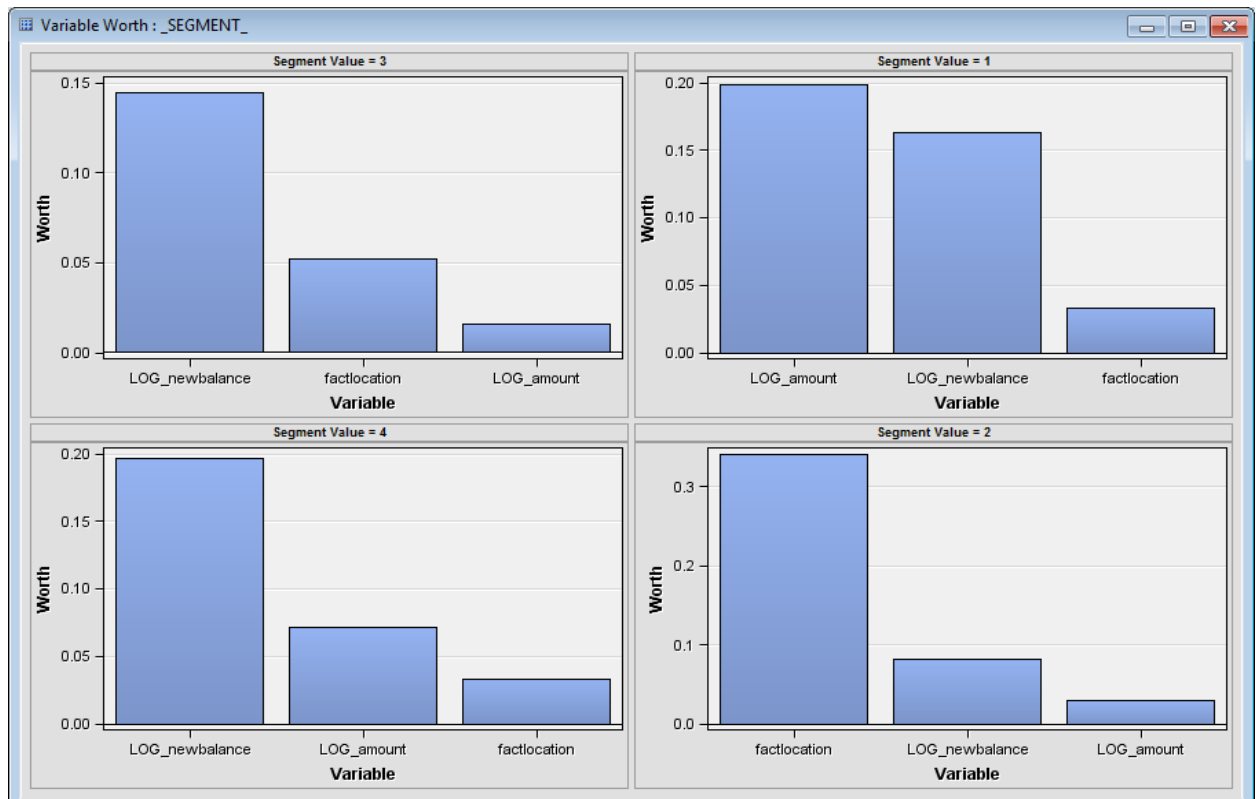


Рисунок 3.17 – Вага кожної змінної у формуванні відповідного кластеру в пакеті SAS Enterprise Miner

Отже, за допомогою кластерного аналізу було досліджено інформацію про проведені транзакції клієнтами мобільного та інтернет-банкінгу. Визначено, що існує певна закономірність між місцеположенням пристроїв, з яких виконувались транзакції, сумами коштів на рахунках клієнтів та балансами після виконання транзакцій. Їх значення та зміни впливають на ознаку втручання у банківську систему.

3.1.2 Побудова регресійних моделей

На рисунку 3.18 представимо загальну діаграму процесу моделювання в пакеті SAS Enterprise Miner та розглянемо нижче побудову моделей регресії, дерева рішень та нейронної мережі на основі повної вибірки вхідних даних.

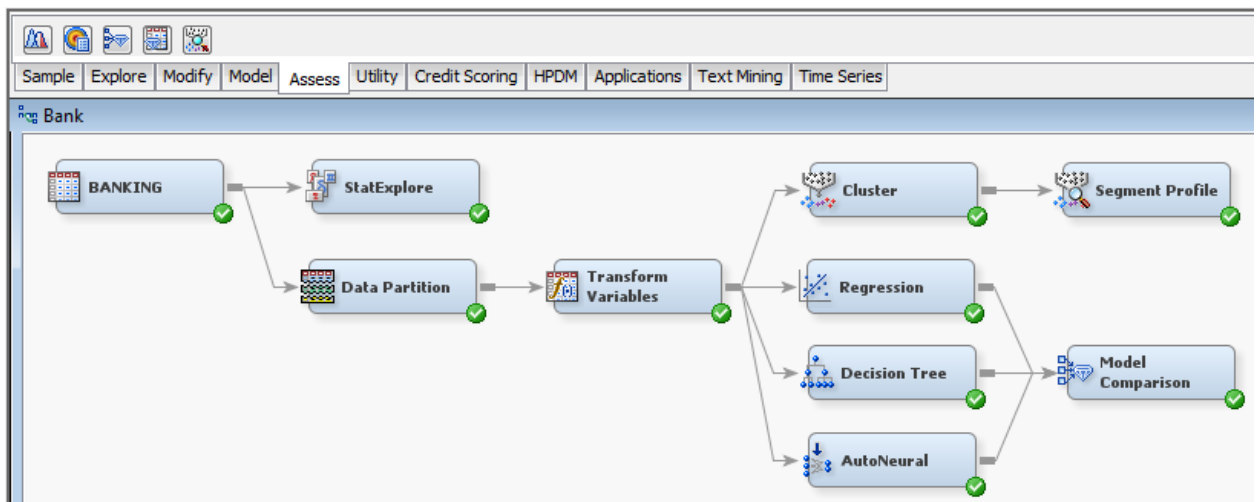


Рисунок 3.18 – Діаграма процесу моделювання в пакеті SAS Enterprise Miner

Додаємо компонент Regression з вкладки Model на робочу область діаграми. З'єднаємо вузол Transform Variables з вузлом Regression. Вузол Regression може створювати кілька типів моделей регресії, включаючи лінійні та логістичні. Тип регресії за замовчуванням визначається рівнем вимірювань цільової змінної.

Виконаємо вузол Regression та переглянемо результати його виконання (рис. 3.19).

Стовпчик таблиці $Pr > ChiSq$ демонструє, наскільки обраний фактор є значущим, а саме: чим менше значення в останньому стовпчику – тим фактор більш значущий у моделі. Змінні з розрахованим значенням < 0.0001 мають високу статистичну значущість, в даному наборі даних це трансформовані змінні балансів до та після проведення транзакцій та змінна місцезположення пристрою, з якого виконувалась транзакція. Інші обрані змінні не є значимими.

Output								
Analysis of Maximum Likelihood Estimates								
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq	Standardized Estimate	Exp (Est)	
Intercept	1	-2.7205	1.8850	2.08	0.1490		0.066	
LOG_amount	1	-0.0193	0.1703	0.01	0.9098	-0.0203	0.981	
LOG_newbalance	1	-0.9849	0.1086	82.29	<.0001	-3.4421	0.373	
LOG_oldbalance	1	0.9038	0.0882	104.98	<.0001	2.8387	2.469	
devicetype I	1	-0.0571	0.1955	0.09	0.7703		0.945	
factlocation Other	1	5.0497	0.2788	328.00	<.0001		155.968	
type CASH_IN	1	1.0186	0.7988	1.63	0.2023		2.769	
type CASH_OUT	1	-0.7451	0.5994	1.55	0.2138		0.475	
type DEBIT	1	-0.0123	1.7023	0.00	0.9942		0.988	
type PAYMENT	1	-0.4049	0.5897	0.47	0.4923		0.667	

Рисунок 3.19 – Проміжний результат виконання вузла Regression

Оберемо метод покрокового виключення незначущих факторів (Stepwise) для побудови моделі логістичної регресії (рис. 3.20).

.. Property	Value
Class Targets	
Regression Type	Logistic Regression
Link Function	Logit
Model Options	
Suppress Intercept	No
Input Coding	Deviation
Model Selection	
Selection Model	Stepwise
Selection Criterion	Default
Use Selection Defaults	Yes
Selection Options	
Optimization Options	

Рисунок 3.20 – Налаштування властивостей вузла Regression

На рисунку 3.21 відображено графік зміни коефіцієнта помилкової класифікації (Misclassification Rate) для навчального та валідаційного набору даних (відповідно синя та червона лінії) в залежності від кроку відбору значущих факторів.

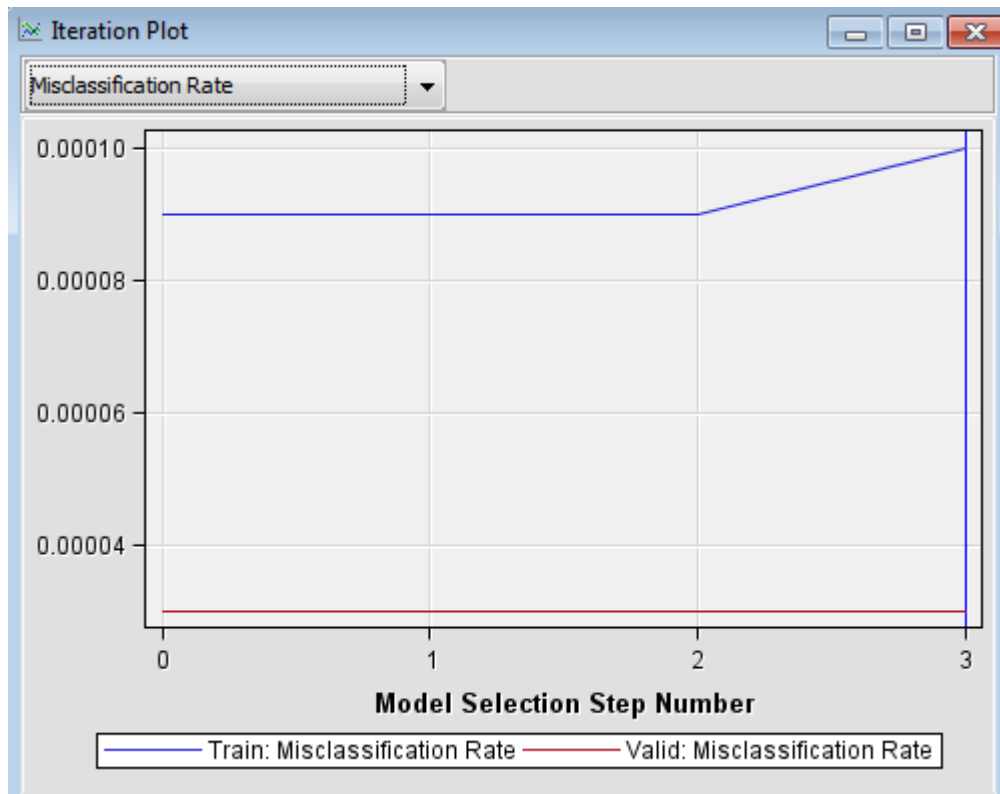


Рисунок 3.21 – Оптимізація логіт-моделі за допомогою коефіцієнта помилкової класифікації в пакеті SAS Enterprise Miner

Результат виконання вузла Regression із значеннями відібраних вхідних змінних зображено на рисунку 2.22.

Output								
Analysis of Maximum Likelihood Estimates								
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq	Standardized Estimate	Exp (Est)	
Intercept	1	-3.4043	0.3518	93.65	<.0001		0.033	
LOG_newbalance	1	-0.8950	0.0910	96.66	<.0001	-3.1280	0.409	
LOG_oldbalance	1	0.8738	0.0846	106.81	<.0001	2.7445	2.396	
factlocation Other	1	5.1102	0.2700	358.11	<.0001		165.707	

Рисунок 3.22 – Результат виконання вузла Regression

Оцінки коефіцієнта ймовірностей, що представлені в стовпчику Exp (Est) результатів виконання вузла Regression, показують, що одинична зміна місцеположення, ініційованого при проведенні транзакції призведе до найбільшої зміни відношення ймовірностей виникнення ознак кіберзагроз.

На рисунку 3.23 представлено графік середньоквадратичної похибки (Average Squared Error) для навчального та валідаційного набору даних.

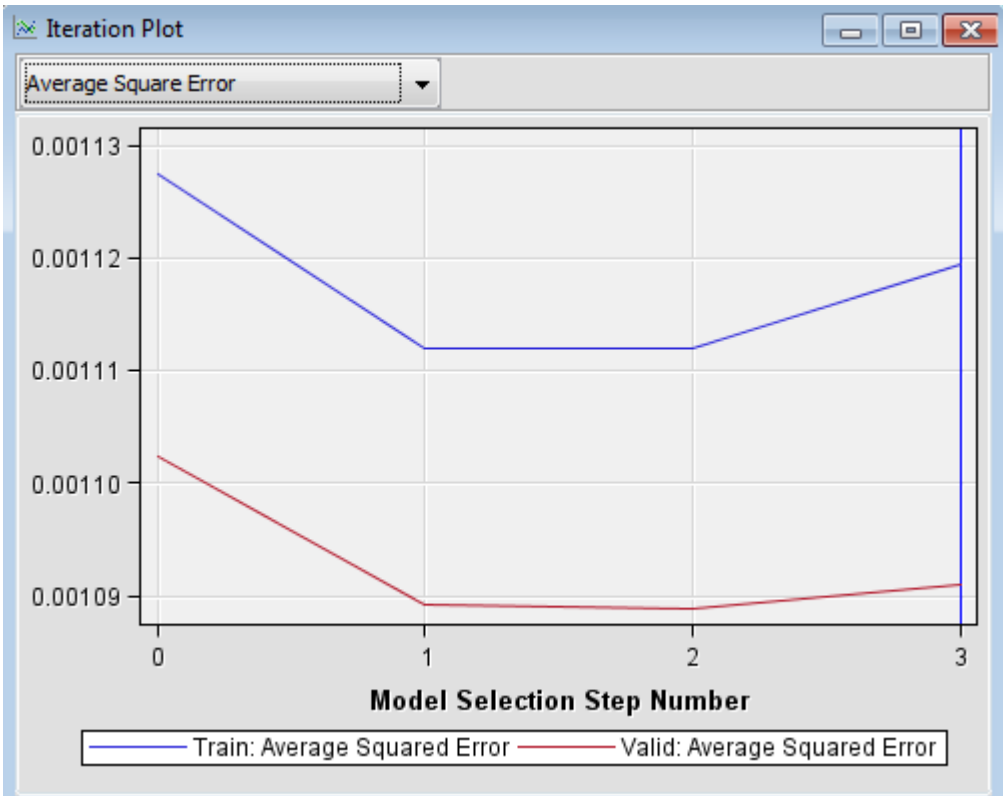


Рисунок 3.23 – Оптимізація логіт-моделі за допомогою середньоквадратичної похибки Average Squared Error (ASE) в пакеті SAS Enterprise Miner

На рисунку 3.24 відображено розраховані значення критерію ксі-квадрата Вальда для обраних факторів (стовпчик Wald Chi-Square).

Output				
500	Type 3 Analysis of Effects			
501				
502			Wald	
503	Effect	DF	Chi-Square	Pr > ChiSq
504				
505	LOG_newbalance	1	96.6623	<.0001
506	LOG_oldbalance	1	106.8145	<.0001
507	factlocation	1	358.1065	<.0001

Рисунок 3.24 – Розрахунок критерію ксі-квадрата Вальда в пакеті SAS Enterprise Miner

Дана модель залишила лише змінні з розрахованим значенням < 0.0001 , що мають високу статистичну значущість.

Таким чином, у результаті покрокового відбору було обрано 3 значущі фактори:

1) Зафіксоване місцеположення пристрою, з якого проводилась транзакція (X_3):

– $X_{3,2}$ – інша країна;

2) Баланс клієнта після проведення транзакції (X_5);

3) Баланс клієнта до проведення транзакції (X_6).

Отже, математична модель вірогідності виникнення кіберзагроз під час проведення транзакцій користувачами мобільного та інтернет-банкінгу має наступний вигляд:

$$\text{logit}(\hat{p}) = -3,4 + 5,11X_{3,2} - 0,89X_5 + 0,87X_6. \quad (3.1)$$

Отже, ймовірність того, що банківська транзакція виявиться кіберзагрозою зростає із присутністю зафіксованого факту проведення транзакції в іншій країні, з великим значенням балансу до проведення транзакції та зменшується із великим значенням балансу після проведення транзакції.

3.1.3 Побудова дерева рішень

Інструмент для побудови дерева рішень в пакеті SAS Enterprise Miner – Decision Tree. Додамо даний вузол в робочу область діаграми, поєднаємо з вузлом Transform Variables та запустимо до виконання, натиснувши Run.

Побудова дерева рішень відбувалась в автоматичному режимі. У результаті чого було згенеровано трирівневе дерево класифікації (рис. 3.25).

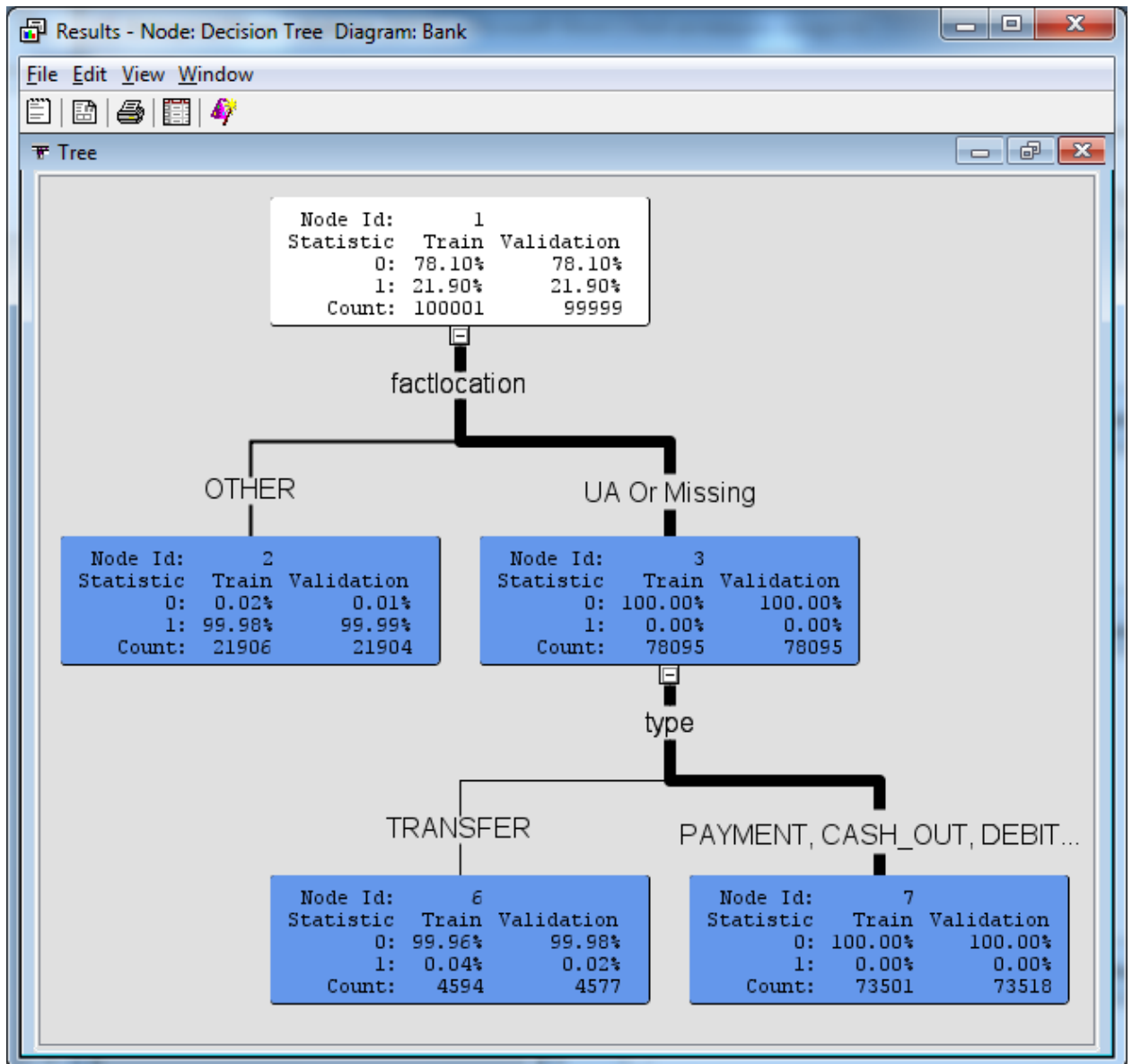


Рисунок 3.25 – Автоматичне дерево рішень в пакеті SAS Enterprise Miner

З побудованої діаграми дерева рішень видно, що найбільш вагомий фактор – місцезположення пристрою, з якого виконувалась транзакція. Після нього за важливістю – тип пристрою, з якого виконувалась транзакція.

На рисунку 3.26 відображено графік зміни коефіцієнта помилкової класифікації (Misclassification Rate) для навчального та валідаційного набору даних (відповідно синя та червона лінії) в залежності від кількості рівнів розгалуження дерева.

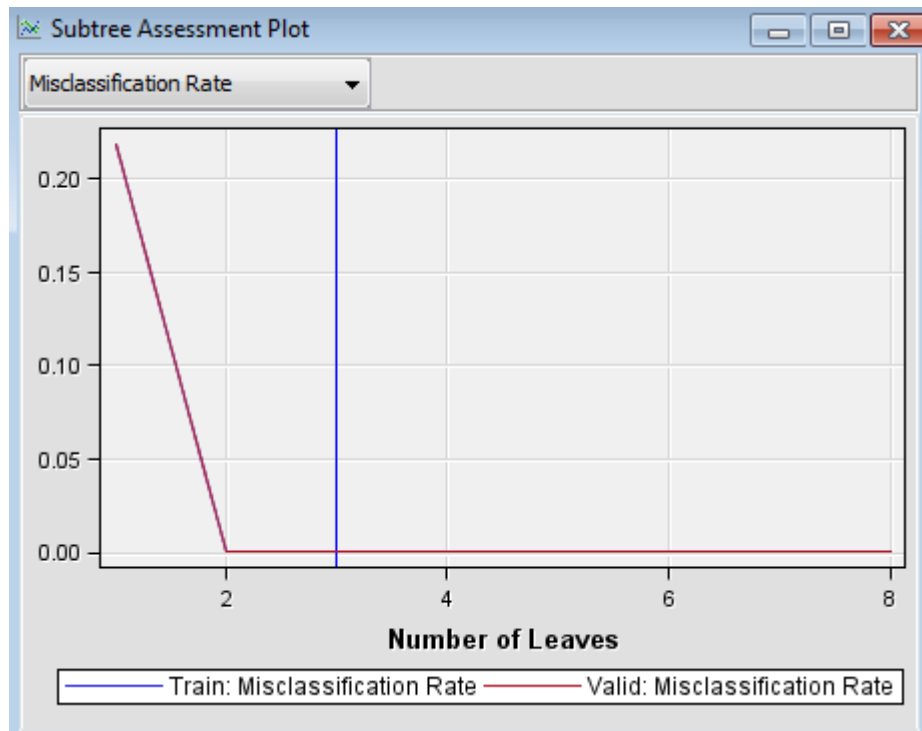


Рисунок 3.26 – Оптимізація автоматичного дерева рішень за допомогою коефіцієнта помилкової класифікації в пакеті SAS Enterprise Miner

Як видно з рисунку 3.26, для тренувального і валідаційного наборів даних графік має спадну тенденцію. Це говорить про те, що зі збільшенням кількості гілок дерева коефіцієнт помилкової класифікації зменшується, тобто, дерево краще класифікує дані. На 3-му кроці цей показник досягає свого мінімального значення для валідаційного набору (це відображено на рис. 3.26 вертикальною синьою лінією), отже подальше нарощування кількості голок не є доцільним. Таким чином, у якості оптимального варіанту було обрано дерево з 3-ма гілками розгалужень.

Таким чином, за результатами побудови дерева рішень, найімовірніше виконана транзакція не є кіберзагрозою, якщо фіксоване місцеположення виконання транзакції клієнтом банкінгу – Україна. А також з'ясовано, що безпечними для користувачів на випадок наявності кіберзагрози є наступні типи загроз: поповнення та зняття коштів, списання коштів з рахунку та проведення оплати.

3.1.4 Побудова моделі на основі нейронної мережі

На діаграму додамо нейронну мережу. Виберемо вкладку Model. Додамо компонент AutoNeural на робочу область діаграми та з'єднаємо вузол Transform Variables з вузлом AutoNeural. Інструмент AutoNeural може використовуватися для автоматичного налаштування нейронної мережі, дозволяє автоматично складати, навчати і перевіряти багатошарові нейронні мережі з прямим зв'язком. Він проводить обмежений пошук кращої конфігурації мережі. У загальному випадку, кожен вхідний елемент повністю пов'язаний з першим прихованим шаром, кожен прихований шар повністю пов'язаний з наступним прихованим шаром, а останній прихований шар повністю пов'язаний з вихідними даними.

Оптимізація мережі відбувалася шляхом мінімізації коефіцієнта помилкової класифікації, на рисунку 3.27 відображено графік його зміни в залежності від кількості нейронів.

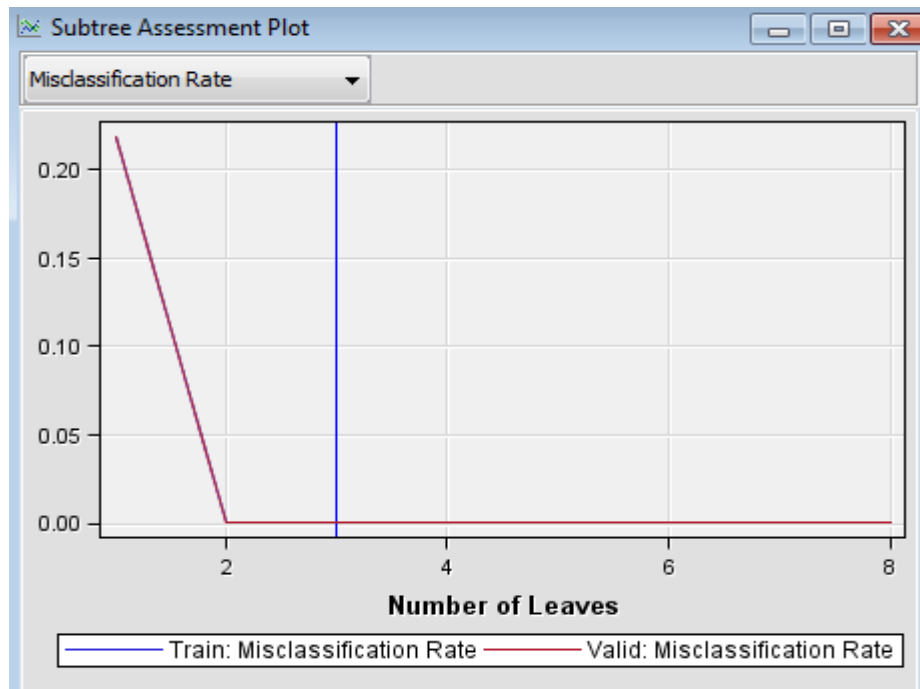


Рисунок 3.27 – Оптимізація нейронної за допомогою коефіцієнта помилкової класифікації в пакеті SAS Enterprise Miner

Як видно, для тренувального і валідаційного наборів графік коефіцієнта стрімко спадає на проміжку (0;2). Потім не змінюються як для навчальних даних так і для валідаційних, а отже побудова мережі з більшою кількістю шарів є недоцільною.

В результаті було згенеровано нейронну мережу, яка складається з 1-го прихованого шару з двома нейронами. Схематично її представлено на рис. 3.28, причому вхідним шаром виступає уся сукупність вхідних даних, а вихідним – значення ймовірності настання кіберзагрози в транзакціях користувачів мобільного та інтернет-банкінгу.

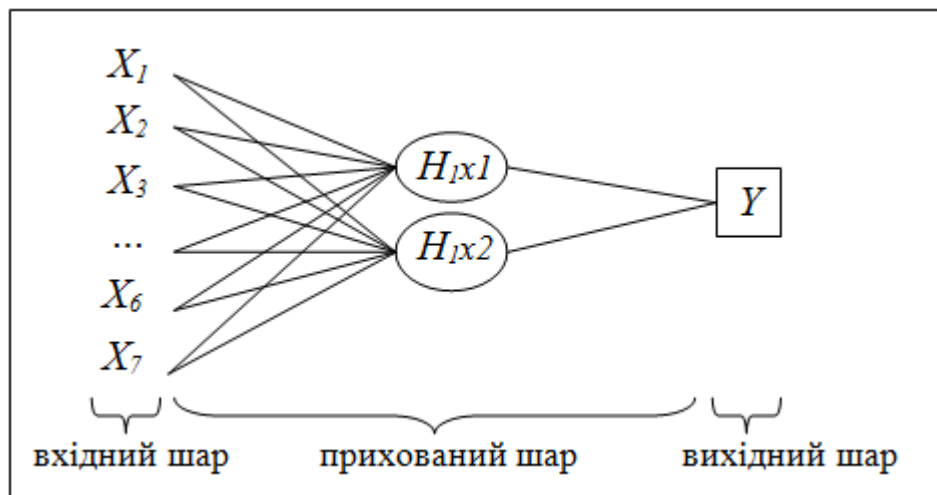


Рисунок 3.28 – Схематичне представлення нейронної мережі

Математична інтерпретація вихідного шару, а також 1-го та 2-го прихованих шарів нейрону наведено у формулах (3.2 – 3.4) відповідно, за допомогою підстановки значень коефіцієнтів змінних у формули (2.5 – 2.6):

$$Y = 1,02 + 7,56 \cdot H_{1x1} - 1,76 \cdot H_{1x2}; \quad (3.2)$$

$$H_1 = \tanh(-0,78 + 0,04 \cdot \text{LOG}X_1 - 0,13 \cdot X_{2,2} + 0,87 \cdot X_{3,2} - 2,8 \cdot \text{LOG}X_5 + 1,36 \cdot \text{LOG}X_6 + 1,78 \cdot X_{7,1} - 0,76 \cdot X_{7,2} - 0,35 \cdot X_{7,3} - 0,58 \cdot X_{7,4}); \quad (3.3)$$

$$H_2 = \tanh(0,99 - 0,05 \cdot LOGX_1 - 0,23 \cdot X_{2,2} - 0,78 \cdot X_{3,2} - 0,11 \cdot LOGX_5 - 0,25 \cdot LOGX_6 + 0,05 \cdot X_{7,1} - 0,07 \cdot X_{7,2} - 0,3 \cdot X_{7,3} - 0,4 \cdot X_{7,4}). \tag{3.4}$$

Вигляд архітектури мережі в пакеті SAS Enterprise Miner наведено на рисунку 3.29.

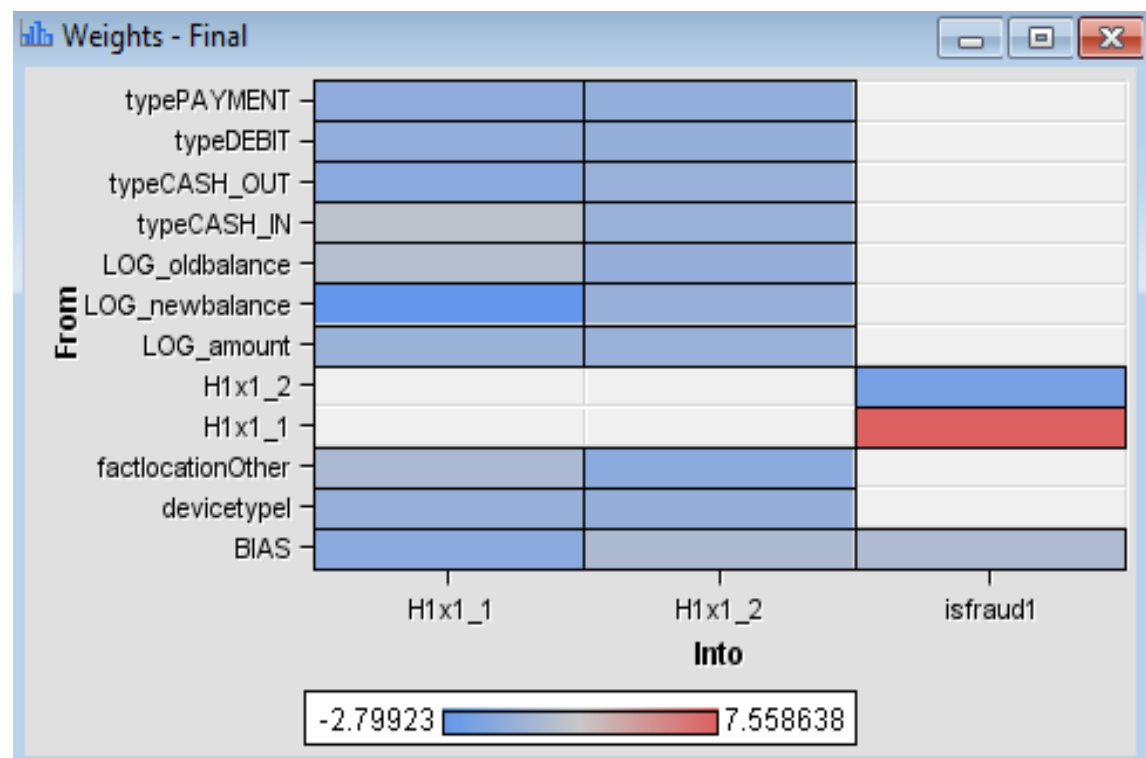


Рисунок 3.29 – Архітектура нейронної мережі в пакеті SAS Enterprise Miner

На рисунку 3.30 наведені вагові оцінки для формул (3.2 – 3.4).

Label	From	Into	Weight
LOG_amount -> H1x1_1	LOG_amount	H1x1_1	0.044417
LOG_newbalance -> H1x1_1	LOG_newbalance	H1x1_1	-2.79923
LOG_oldbalance -> H1x1_1	LOG_oldbalance	H1x1_1	1.360785
LOG_amount -> H1x1_2	LOG_amount	H1x1_2	-0.05355
LOG_newbalance -> H1x1_2	LOG_newbalance	H1x1_2	-0.11025
LOG_oldbalance -> H1x1_2	LOG_oldbalance	H1x1_2	-0.25139
devicetypel -> H1x1_1	devicetypel	H1x1_1	-0.13465
factlocationOther -> H1x1_1	factlocationOther	H1x1_1	0.867504
typeCASH_IN -> H1x1_1	typeCASH_IN	H1x1_1	1.775061
typeCASH_OUT -> H1x1_1	typeCASH_OUT	H1x1_1	-0.75885
typeDEBIT -> H1x1_1	typeDEBIT	H1x1_1	-0.34715
typePAYMENT -> H1x1_1	typePAYMENT	H1x1_1	-0.57974
devicetypel -> H1x1_2	devicetypel	H1x1_2	-0.23464
factlocationOther -> H1x1_2	factlocationOther	H1x1_2	-0.78262
typeCASH_IN -> H1x1_2	typeCASH_IN	H1x1_2	0.048199
typeCASH_OUT -> H1x1_2	typeCASH_OUT	H1x1_2	-0.0721
typeDEBIT -> H1x1_2	typeDEBIT	H1x1_2	-0.30449
typePAYMENT -> H1x1_2	typePAYMENT	H1x1_2	-0.39577
BIAS -> H1x1_1	BIAS	H1x1_1	-0.77711
BIAS -> H1x1_2	BIAS	H1x1_2	0.991864
H1x1_1 -> isfraud1	H1x1_1	isfraud1	7.558638
H1x1_2 -> isfraud1	H1x1_2	isfraud1	-1.75976
BIAS -> isfraud1	BIAS	isfraud1	1.022777

Рисунок 3.30 – Вагові коефіцієнти нейронної мережі

Нейронні мережі мають властивість адаптовуватися до змін навколишнього середовища, тобто в нестационарних умовах, коли інформація змінюється з часом. Цю властивість як раз доцільно використовувати у випадку створення нейронної мережі для аналізу банківських транзакцій, зміни в яких відбуваються постійно.

3.2 Аналіз якості та адекватності побудованих моделей

Аналіз якості моделі, яка прогнозує ймовірність настання певної події, визначається, в першу чергу, за тим, наскільки добре вона передбачила результат. Такі характеристики визначаються кількісно, у відсотковому відношенні, а також за коефіцієнтом помилкової класифікації (Misclassification Rate, *MISC*), який визначається за формулою (3.5):

$$MISC = \frac{Nm}{N}, \quad (3.5)$$

де Nm - кількість неправильно класифікованих випадків;
 N - загальна кількість випадків.

Показником точності моделі є середньоквадратична похибка (Mean Squared Error, MSE), яка розраховується за формулою (3.6).

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}, \quad (3.6)$$

де \hat{y}_i – змодельоване значення (прогнозована ймовірність настання події);
 y_i – фактичне значення показника;
 n – порядковий номер [31].

Доцільним є дослідження результатів класифікації саме для валідаційних даних, адже на тренувальних модель будується, тому якісний результат заздалегідь є очевидним, а на валідаційних даних перевіряються прогностичні властивості побудованої моделі.

В таблиці 3.3 представлено характеристики класифікації за допомогою регресійної моделі на навчальній вибірці.

Таблиця 3.3 – Характеристика класифікаційних властивостей регресійної моделі на навчальній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9923	99,9949	78093	78,0922
1	0	0,0077	0,0274	6	0,0060
0	1	0,0183	0,0051	4	0,0040
1	1	99,9817	99,9726	21898	21,8978

Отже, з таблиці бачимо, що модель на навчальній вибірці вірно класифіковано 99,99% транзакцій, які не виявились кіберзагрозами і 99,97% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,03% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,01% транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78093 транзакціях, а позитивний – 21898. Неправильно класифіковано негативний результат у 4 транзакціях, а позитивний – у 6. У цілому можна сказати, що частка правильної класифікації склала 99,99% (78,0922% + 21,8978%).

Таблиця 3.4 демонструє характеристики класифікації регресійної моделі на валідаційних даних. Остання вірно класифікувала 99,99% транзакцій, які не виявились кіберзагрозами і 99,99% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,005% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,002% транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Таблиця 3.4 – Характеристика класифікаційних властивостей регресійної моделі на валідаційній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9987	99,9974	78094	78,0948
1	0	0,0013	0,0046	1	0,0010
0	1	0,0091	0,0026	2	0,0020
1	1	99,9909	99,9954	21902	21,9022

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78094 транзакціях, а позитивний – 21902. Неправильно класифіковано негативний результат у 2 транзакціях, а позитивний – у 1. У цілому можна сказати, що частка правильної класифікації склала 99,997% (78,0948% + 21,9022%).

В таблиці 3.5 представлені основні коефіцієнти, що характеризують якість регресійної моделі.

Таблиця 3.5 – Коефіцієнти якості регресійної моделі

Коефіцієнт	Вибірка	
	Навчальна	Валідаційна
Частка неправильної класифікації (Misclassification Rate, MISC)	0,0001	0,00003
Середньоквадратична похибка (Mean Square Error, MSE)	0,001119	0,001091
Середньоквадратична похибка (Average Squared Error, ASE)	0,001119	0,001091

Таким чином, досить низькі значення розрахованих коефіцієнтів свідчать про якість та адекватність побудованої моделі.

В таблиці 3.6 наведено характеристики класифікації на основі дерева рішень на навчальній вибірці.

Таблиця 3.6 – Характеристика класифікаційних властивостей дерева рішень на навчальній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9974	99,9949	78093	78,0922
1	0	0,0026	0,0091	2	0,0020
0	1	0,0183	0,0051	4	0,0040
1	1	99,9817	99,9909	21902	21,9018

Отже, з таблиці бачимо, що модель на навчальній вибірці вірно класифіковано 99,99% транзакцій, які не виявились кіберзагрозами і 99,99% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,009% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,005% транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78093 транзакціях, а позитивний

– 21902. Неправильно класифіковано негативний результат у 4 транзакціях, а позитивний – у 2. У цілому можна сказати, що частка правильної класифікації склала 99,994% (78,0922% + 21,9018%).

Таблиця 3.7 демонструє характеристики класифікації на основі дерева рішення на валідаційних даних. Остання вірно класифікувала 99,99% транзакцій, які не виявились кіберзагрозами і 99,99% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,005% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,002% транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Таблиця 3.7 – Характеристика класифікаційних властивостей дерева рішень на валідаційній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9987	99,9974	78094	78,0948
1	0	0,0013	0,0046	1	0,0010
0	1	0,0091	0,0026	2	0,0020
1	1	99,9909	99,9954	21902	21,9022

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78094 транзакціях, а позитивний – 21902. Неправильно класифіковано негативний результат у 2 транзакціях, а позитивний – у 1. У цілому можна сказати, що частка правильної класифікації склала 99,997% (78,0948% + 21,9022%).

В таблиці 3.8 подано коефіцієнти, що описують якість дерева рішень.

Таблиця 3.8 – Коефіцієнти якості дерева рішень

Коефіцієнт	Вибірка	
	Навчальна	Валідаційна
Частка неправильної класифікації (Misclassification Rate, MISC)	0,00009	0,00003
Середньоквадратична похибка (Average Squared Error, ASE)	0,001112	0,001097

Таким чином, можна зробити висновок, що побудоване дерево рішень є якісним й адекватним, що підтверджується аналізом результатів моделювання.

Таблиці 3.9 демонструє характеристики класифікації на основі нейронної мережі на навчальній вибірці.

Таблиця 3.9 – Характеристика класифікаційних властивостей нейронної мережі на навчальній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9949	99,9987	78096	78,0952
1	0	0,0051	0,0183	4	0,0040
0	1	0,0046	0,0013	1	0,0010
1	1	99,9954	99,9817	21900	21,8998

Отже, з таблиці бачимо, що модель на навчальній вибірці вірно класифіковано 99,99% транзакцій, які не виявились кіберзагрозами і 99,98% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,018% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,001% транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78096 транзакціях, а позитивний – 21900. Неправильно класифіковано негативний результат у 1 транзакції, а позитивний – у 4. У цілому можна сказати, що частка правильної класифікації склала 99,995% (78,0952% + 21,8998%).

Таблиця 3.10 демонструє характеристики класифікації на основі дерева рішення на валідаційних даних. Остання вірно класифікувала 99,99% транзакцій, які не виявились кіберзагрозами і 99,99% транзакцій, які є кіберзагрозами. У цей же час модель класифікувала 0,005% транзакцій, що були кіберзагрозами, як ті, що не виявились кіберзагрозами і 0,001%

транзакцій, які не виявились кіберзагрозами, було класифіковано, як ті, що є кіберзагрозами.

Таблиця 3.10 – Характеристика класифікаційних властивостей нейронної мережі на валідаційній вибірці

Цільова змінна	Результат	Цільова змінна, %	Результат, %	Частота випадків	Загальна класифікація, %
0	0	99,9987	99,9987	78095	78,0958
1	0	0,0013	0,0046	1	0,0010
0	1	0,0046	0,0013	1	0,0010
1	1	99,9954	99,9954	21902	21,9022

Щодо абсолютних величин, то модель правильно класифікувала негативний результат (не є кіберзагрозою) у 78095 транзакціях, а позитивний – 21902. Неправильно класифіковано негативний результат у 1 транзакції та позитивний – у 1. У цілому можна сказати, що частка правильної класифікації склала 99,998% (78,0958% + 21,9022%).

В таблиці 3.11 представлені основні коефіцієнти, що характеризують якість дерева рішень.

Таблиця 3.11 – Коефіцієнти якості нейронної мережі

Коефіцієнт	Вибірка	
	Навчальна	Валідаційна
Частка неправильної класифікації (Misclassification Rate, MISC)	0,00005	0,00002
Середньоквадратична похибка (Mean Square Error, MSE)	0,001105	0,001094
Середньоквадратична похибка (Average Squared Error, ASE)	0,001105	0,001094

Таким чином, можна зробити висновок, що побудована нейронна мережа є якісною й адекватною, що підтверджується чисельними характеристиками результатів моделювання.

Для порівняння роботи всіх побудованих моделей додамо вузол Model Comparison та з'єднаємо його з вузлами Regression, Decision Tree та AutoNeural (рис. 3.31).

Запустимо на виконання вузол Model Comparison. На цьому етапі дослідження здійснюється порівняльний аналіз побудованих моделей та вибір кращої з них. Відбір найбільш якісної та точної моделі проводиться на основі мінімізації частки неправильної класифікації (3.4) та середньоквадратичної похибки (3.5).

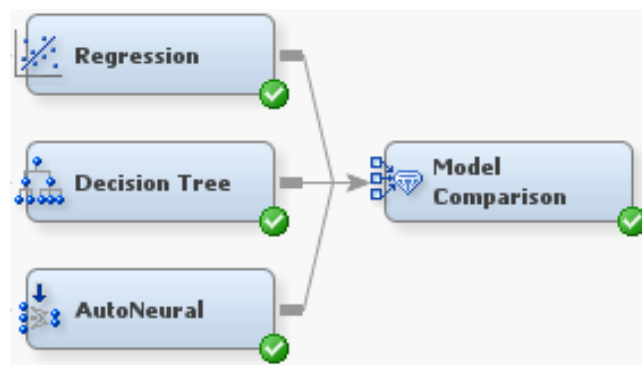


Рисунок 3.31 – Додавання вузла Model Comparison на діаграму

Порівняльна характеристика коефіцієнтів якості та точності моделей: регресійної, дерева рішень та нейронної мережі представлені в таблиці 3.12.

Таблиця 3.12 – Порівняльна характеристика моделей

№ з/п	Модель	Частка неправильної класифікації (Misclassification Rate, MISC)		Середньоквадратична похибка (Mean Square Error, MSE)	
		Валідаційна	Навчальна	Валідаційна	Навчальна
1	AutoNeural	0,00002	0,00005	0,001094	0,001105
2	Decision Tree	0,00003	0,00009	0,001097	0,001112
3	Regression	0,00003	0,0001	0,001091	0,001119

Моделі, представлені в таблиці 3.12 розташовані від найкращої до найгіршої за кількісними оцінками частки неправильної класифікації та

середньоквадратичної похибки. Модель тим краще описує набір даних, чим менші значення цих показників.

Найнижчими значеннями коефіцієнтів частки неправильної класифікації та середньої квадратичної помилки та найкращою моделлю виявилась AutoNeural, на другому місці – Decision Tree та на третьому – Regression.

ROC-крива (Receiver Operating Characteristic, ROC) – графічна характеристика якості бінарного класифікатора, яка відображає залежність частки вірних позитивних класифікацій від частки помилкових позитивних класифікацій при варіюванні порога вирішального правила.

Результат розрахованих значень коефіцієнтів підкріплюється графіками ROC-кривих (рис. 3.32). На рисунку відображено криві для 2-х моделей: навчального та валідаційного наборів даних. Синьою лінією зображено криву дерева рішень, червоною – регресії, а зеленою – нейронної мережі. Чим крива більше віддаляється від базової лінії, тим краще модель класифікує дані, тобто прогнозує можливість виникнення кіберзагрози.

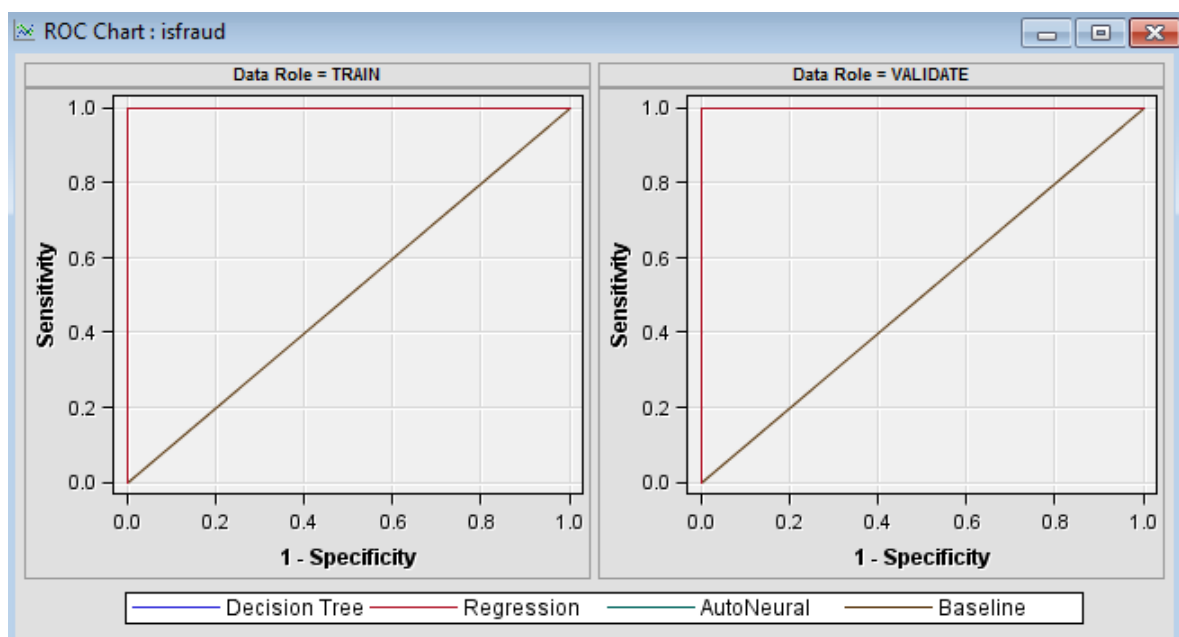


Рисунок 3.32 – ROC-криві дерева рішень, регресії та нейронної мережі для навчального та валідаційного наборів даних у пакеті SAS Enterprise Miner

За даними кривими видно, що ROC-криві моделей накладаються одна на одну, це свідчить про приблизно однакову якість класифікації моделей.

Таким чином, за допомогою інструмента Model Comparison, було визначено, що нейронна мережа краще моделює оцінку результату виникнення кіберзагрози під час проведення банківської транзакції клієнтами мобільного та інтернет-банкінгу, ніж регресія й дерево рішень.

3.3 Оцінка результатів та ефекту

Здійснимо прогнозування за обраною моделлю. Для цього підключимо нове джерело вхідних даних.

Однією з переваг пакету SAS Enterprise Miner в процесі реалізації моделі є засіб скорингу, який здатен додавати прогнози до будь-якого набору даних, структурованого так само, як і навчальні дані. SAS Enterprise Miner пропонує застосування моделі за допомогою сукупності даних з внутрішнім скорингом. Для цього потрібно задати джерело скорингових даних, інтегрувати джерело скорингових даних з інструментом Score в діаграмі ходу процесу і перенести набір скорингових даних в бібліотеку.

Створимо джерело скорингових даних `testbanking.sas7bdat`, виконавши аналогічні дії до створення джерела даних для моделі.

File > New > Data Source > Next > Browse > `testbanking.sas7bdat` > Next.

На кроці 4 Metadata Advisor Options натиснемо Advanced > Customize > змінимо значення властивостей, аналогічно до зміни значень властивостей для моделі.

На кроці 6 Data Source Wizard оберемо роль: Role > Score (рис. 3.33)

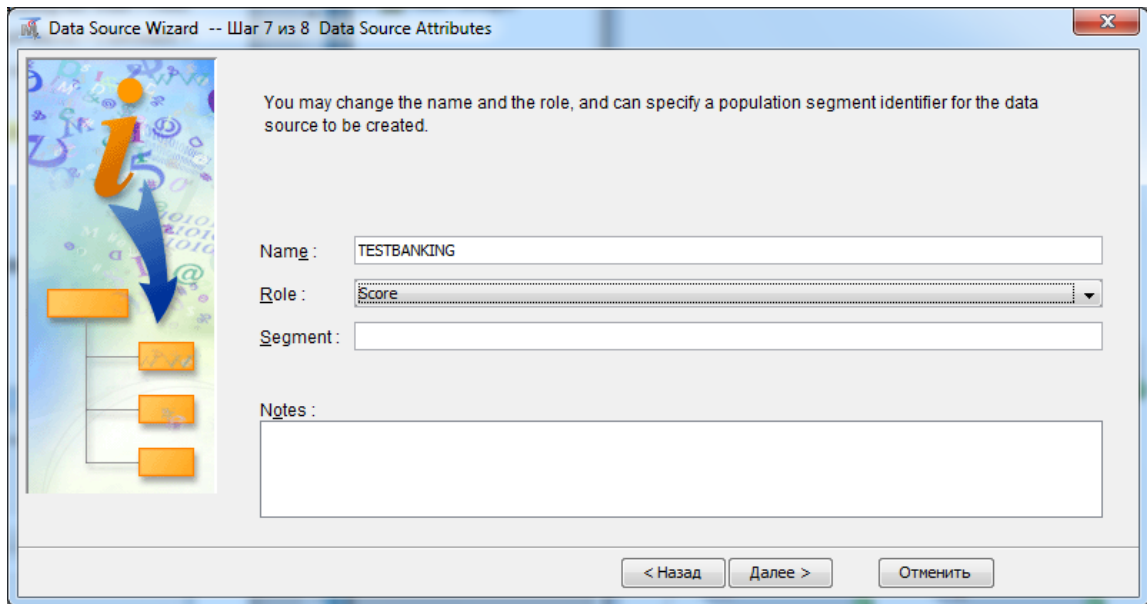


Рисунок 3.33 – Вибір ролі джерела даних

Додамо в робочу область діаграми інструмент Score із вкладки Assess, який приєднує прогнози обраної моделі до набору скорингових даних та зєднаємо його з вузлом моделі, що була обрана найкращою – AutoNeural (рис. 3.34).



Рисунок 3.34 – Додавання вузла Score

В робочу область діаграми перетягнемо джерело даних TESTBANKING та зєднаємо його з вузлом Score (рис. 3.35)

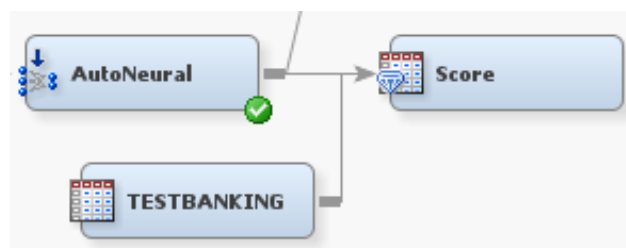
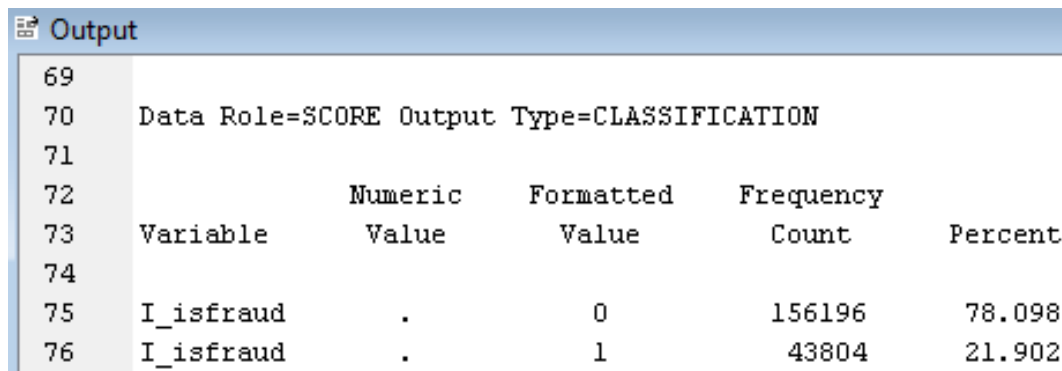


Рисунок 3.35 – Додавання джерела даних TESTBANKING

Запустимо вузол Score та переглянемо результати. У вікні Output представлено частоту рішень для кожного набору даних, який проглядається вузлом Score (рис. 3.36).

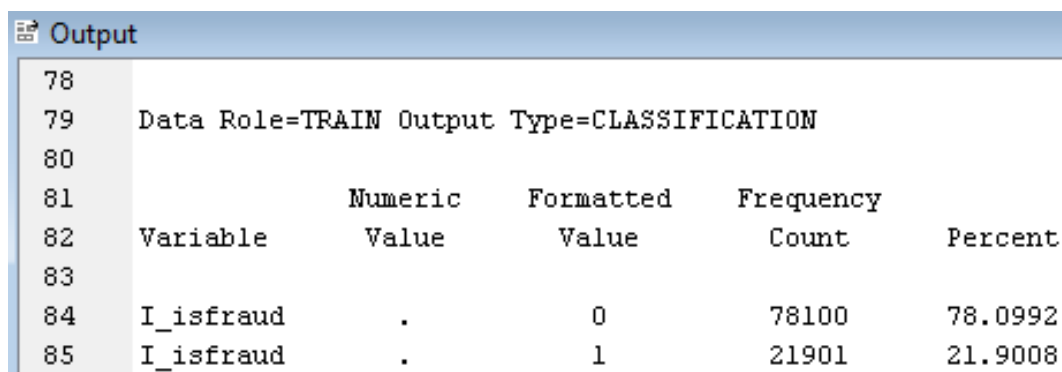


Variable	Numeric Value	Formatted Value	Frequency Count	Percent
I_isfraud	.	0	156196	78.098
I_isfraud	.	1	43804	21.902

Рисунок 3.36 – Процентна частка спостережень скорингового набору

Отже, в скоринговому наборі даних 21,902% транзакцій ймовірно є кіберзагрозами та 78,098% - ні.

На рисунку 3.37 представлено процентну частку спостережень навчального набору даних ймовірності виникнення кіберзагрози.



Variable	Numeric Value	Formatted Value	Frequency Count	Percent
I_isfraud	.	0	78100	78.0992
I_isfraud	.	1	21901	21.9008

Рисунок 3.37 – Процентна частка спостережень навчального набору

В даному наборі кіберзагрозами виявились 21,9008% транзакцій, а 78,0992% – ні.

На рисунку 3.38 представлено процентну частку спостережень валідаційного набору даних ймовірності виникнення кіберзагрози.

Output					
87					
88	Data Role=VALIDATE Output Type=CLASSIFICATION				
89					
90		Numeric	Formatted	Frequency	
91	Variable	Value	Value	Count	Percent
92					
93	I_isfraud	.	0	78096	78.0968
94	I_isfraud	.	1	21903	21.9032

Рисунок 3.38 – Процентна частка спостережень валідаційного набору

У валідаційному наборі даних 21,9032% транзакцій виявились кіберзагрозами, а 78,0968% – ні.

Відсоткові співвідношення майже рівні – це свідчить про стаціонарність між навчальними, валідаційними та скоринговими даними.

Побудована модель виявлення кібернетичних загроз може бути впроваджена в банківських установах з метою попередження потенційних кіберзагроз.

Впровадження отриманої моделі передбачає отримання банком соціального та економічного ефектів.

Соціальний ефект полягає у:

- зростанні рівня довіри клієнтів до банків через підвищення захищеності та надійності;
- можливості вибору клієнтами тих банків, які пропонують послуги, захищені від кіберзагроз.

Економічний ефект передбачає:

- зниження витрат банків від кібернетичних загроз у зв'язку з їх попередженням;

– збільшення прибутку від залучення нових клієнтів за рахунок зростання рівня довіри до банків.

Спрогнозуємо ймовірний обсяг доходів та витрат банку від впровадження моделі виявлення ознак кіберзагроз.

За 2017 рік банк, дані якого було досліджено, мав збитки у розмірі 240 тис. грн.

Розрахуємо можливий розмір доходу банку від впровадження моделі. Банком буде отримано дохід, за рахунок скорочення витрат на усунення наслідків кібератак на 1%, в розмірі (3.7):

$$D_1 = 240000 \cdot 0,01 = 2400 \text{ (грн.)} \quad (3.7)$$

Доходи, отримані банком в 2017 році у результаті здійснення операцій чи надання послуг своїм клієнтам, склали 48 млн. грн.

Розрахуємо можливий розмір доходу банку від збільшення кількості залучення нових клієнтів на 1% (3.8):

$$D_2 = 48000000 \cdot 0,01 = 480000 \text{ (грн.)} \quad (3.8)$$

Отже, загальна сума річної економії (S) банку становитиме (3.9):

$$S = D_1 + D_2 = 2400 + 480000 = 482400 \text{ (грн.)} \quad (3.9)$$

Проведемо розрахунки суми капітальних витрат (C) на впровадження моделі виявлення ознак кібернетичних загроз, якщо витрати на заробітну плату розробників становлять 310 тис. грн., амортизаційні витрати – 75 тис. грн., витрати на обслуговування – 190 тис. грн. (3.10):

$$C = 310000 + 75000 + 190000 = 575000 \text{ (грн.)} \quad (3.10)$$

Таким чином, можемо розрахувати річний економічний ефект (E_y), за умови що нормативний коефіцієнт окупності капітальних вкладень (m) становить 0,33 (3.11):

$$E_y = S - C \cdot m = 482400 - 575000 \cdot 0,33 = 292650 \text{ (грн.)} \quad (3.11)$$

Визначимо коефіцієнт ефективності (R_{ce}) капітальних вкладень (3.12):

$$R_{ce} = \frac{S}{C} = \frac{482400}{575000} = 0,84 \quad (3.12)$$

Визначимо термін окупності (To) витрат на впровадження моделі (3.13):

$$To = \frac{1}{R_{ce}} = \frac{C}{S} = \frac{575000}{482400} = 1,2 \text{ (років)} \quad (3.13)$$

Отже, аналіз ефективності моделі для банку показав, що річний економічний ефект становитиме 292650 грн., а термін окупності витрат на впровадження моделі – 1 рік та 72 дні.

Показники ефективності моделі є досить високими, що свідчить про можливість її практичного застосування в діяльності банківських установ.

ВИСНОВКИ

У роботі було розкрито сутність обраного для моделювання об'єкту – взаємовідносин учасників банківської діяльності, в результаті яких створюються умови для виникнення кіберзагроз, що становлять небезпеку банківській сфері. Виявлено основні групи кіберзагрози в банках: атаки мережевого та прикладного рівнів; соціальна інженерія; розвинені стійкі загрози; організована кіберзлочинність, порушення основних даних. Значної шкоди банки зазнають від фішингових атак та найбільш поширеною мобільною кіберзагрозою є банківські трояни, оскільки в більшості володарів смартфонів є в наявності і банківська карта.

Проаналізовано існуючі підходи до виявлення кіберзагроз та з'ясовано, що основними заходами українських банків є збільшення штату співробітників служби безпеки, залучення додаткових ресурсів для проведення розслідувань. В той час, як іноземні банки застосовують більш інноваційні підходи стосовно даного питання, запроваджуючи інструменти бізнес-аналітики. Одним з таких підходів є гібридний, сутність якого полягає в комбінуванні різних методів і алгоритмів: експертних і статистичних бізнес-правил, моделі відхилення від звичайної схеми поведінки клієнта, пошуку прихованих закономірностей в даних, аналізу соціальних мереж.

Досліджено методи інтелектуального аналізу, серед яких виділяють три групи: технологічні, статистичні та кібернетичні. Найбільш перспективним напрямком інтелектуального аналізу є кібернетичні методи, які представляють собою множину підходів, об'єднаних ідеєю комп'ютерної математики та використання теорії штучного інтелекту. Серед цієї групи методів в ході аналізу було обрано для побудови моделі: регресію, дерево прийняття рішень та нейронну мережу.

В процесі підготовки до побудови моделі виявлення ознак кіберзагроз у банку в якості вихідних даних було використано інформацію, що міститься у базі даних мобільного та інтернет-банкінгу банку «Х». Оскільки дана інформація є комерційною таємницею, то розголошення назви банківської установи не є можливим.

Основою для моделі було обрано вісім вхідних змінних, серед яких одна – цільова, виражена випадками виявлення ознак кіберзагроз.

Первинний та кластерний аналіз вхідних даних, побудова моделей відбувались з використанням програмного пакету SAS Enterprise Miner.

В результаті проведеного первинного аналізу було отримано основні статистичні характеристики вхідних змінних, визначено ролі змінних у моделюванні, а також виявлено, що у вхідному масиві даних відсутні пропущені значення в інтервальних змінних. Було проаналізовано графіки інтервальних змінних, які показали, що розподіл даних величин не відповідає нормальному закону розподілу, у зв'язку з чим інтервальні вхідні змінні було прологорифмовано.

Для виявлення прихованих, неочевидних тенденцій та закономірностей у вхідних даних було проведено більш серйозний, глибинний статистичний аналіз – кластерний, результатом якого було виділення чотирьох кластерів. Профільне дослідження кластерів показало, що при формуванні першого кластеру найбільшу вагу мали змінні, що містять баланс клієнта до та після проведення транзакції. На формування другого кластеру найбільше вплинула змінна, яка відображає ініційоване місцеположення пристрою, з якого проведено транзакцію. Змінна, яка містить суму, що знаходиться на балансі клієнта після проведення транзакції спричинила значний вплив на формування третього та четвертого кластерів.

Для побудови регресійної моделі, а саме логістичної регресії, було обрано метод покрокового виключення незначущих факторів (Stepwise). У результаті покрокового відбору було обрано 3 значущі фактори: зафіксоване

місцеположення пристрою, з якого проводилась транзакція; баланс клієнта після проведення транзакції та баланс клієнта до проведення транзакції.

Побудова дерева рішень відбувалась в автоматичному режимі. У результаті чого було згенеровано трирівневе дерево класифікації, яке вказало на найбільш вагомні фактори – місцеположення пристрою, з якого виконувалась транзакція та тип пристрою, з якого виконувалась транзакція.

Було згенеровано нейронну мережу, яка складається з 1-го прихованого шару з двома нейронами та представлено математичну інтерпретацію вихідного шару, а також 1-го та 2-го прихованих шарів нейрону.

Побудовані моделі пройшли перевірку на адекватність та якість з використанням коефіцієнтів: частки неправильної класифікації MISC, середньоквадратичних похибок MSE та ASE. Проаналізовано їх результати та встановлено, що усі побудовані моделі майже однаково точно описують вхідні дані, проте за усіма показниками найкращою виявилась модель на основі нейронної мережі.

Було проведено прогнозування ймовірності виникнення ознак кіберзагрози під час проведення транзакції користувачами мобільного та інтернет-банкінгу на основі обраної моделі нейронної мережі. В результаті чого, виявилось, що 21,9% транзакцій ймовірно містить ознаки кіберзагрозами.

Застосування данної моделі допоможе працівникам банківського сектору виявляти в транзакціях ознаки кібернетичних загроз, тим самим попереджуючи користувачів мобільного та інтернет-банкінгу від можливих збитків, завданих злочинними діями.

Проте дана модель потребує постійного оновлення та удосконалення у зв'язку з появою нових загроз для користувачів мобільного та інтернет-банкінгу. Необхідно доповнювати вибірку даних актуальною інформацією про виконанні користувачами транзакції. Використання моделі в банках призведе до отримання соціального та економічного ефектів.

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Anil K. Jain with RadhaChitta and Rong Jin, Clustering Big Data: Department of Computer Science Michigan State University, 2012
2. Cerrito PB: Introduction to Data Mining Using SAS Enterprise Miner, SAS Institute Inc., 2008.
3. Denial of Service Attacks [Електронний ресурс] – Режим доступу: <https://s2.ist.psu.edu/paper/DDoS-Chap-Gu-June-07.pdf>
4. DoS-атака [Електронний ресурс] – Режим доступу: <https://ru.wikipedia.org/wiki/DoS-атака>
5. IT threat evolution Q3 2017. Statistics [Електронний ресурс]. – Режим доступу : <https://securelist.com/it-threat-evolution-q3-2017-statistics/83131/>
6. Jiawei Han, Micheline Kamber, Jian Pei. Data Mining: Concepts and Techniques, 2012
7. K. A. Abdul Nazeer, M. P. Sebastian. Improving the Accuracy and Efficiency of the k-means Clustering Algorithm, London, U.K, 2009
8. K. Munivara Prasad, A. Rama Mohan Reddy, K. Venugopal Rao. DoS and DDoS Attacks: Defense, Detection and Traceback Mechanisms – A Survey [Електронний ресурс] – Режим доступу: https://globaljournals.org/GJCST_Volume14/3-DoS-and-DDoS-Attacks-Defense-Detection.pdf
9. Madan Lal Bhasin. Data Mining:A Competitive Tool in the Banking and Retail Industries, The Chartered Accountant October, 2006.
10. SAS Enterprise Miner. Обзор решения [Електронний ресурс] – Режим доступу: https://www.sas.com/content/dam/SAS/ru_ru/doc/factsheet/sas-enterprise-miner-04-04-2016.pdf

11. Simon Haykin. Neural Networks and Learning Machines Third Edition – University Hamilton, Ontario, Canada, 2008. – 906 с.
12. The Top Five Security Threats to Your Banking Institution [Електронний ресурс]. – Режим доступу : http://www.level3.com/-/media/files/infographics/en_infg_financialserv_topnetworksecuritythreats_regionalbanks.pdf
13. Trend Report «Financial Cyber Threats Q1 2017» conducted with Kaspersky Labs and Telefónica [Електронний ресурс]. – Режим доступу : http://www.level3.com//media/files/infographics/en_infg_financialserv_topnetworksecuritythreats_regionalbanks.pdf
14. Барсегян А. А., Куприянов М. С., Степаненко В. В., Холод И. И. Б26. Методы и модели анализа данных: OLAP и Data Mining. — СПб.: БХВ-Петербург, 2004. — 336 с.
15. Бахрушин В.Є. Методи аналізу даних : навчальний посібник для студентів /. В.Є. Бахрушин. – Запоріжжя : КПУ, 2011. – 268 с.
16. Введение в анализ ассоциативных правил [Електронний ресурс]. – Режим доступу : <https://basegroup.ru/community/articles/intro>
17. Гибридный подход к обнаружению и предотвращению мошенничества [Електронний ресурс]. – Режим доступу : <http://channel4it.com/persons/Gibridnyy-pohod-k-obnaruzheniyu-i-predotvrashcheniyu-moshennichestva-v-bankah-8574.html#>
18. Головна мобільна кіберзагроза [Електронний ресурс]. – Режим доступу : <http://www.ohrana-ua.com/articles/837-golovna-mobl-na-kberzagroza.html>
19. Дерево прийняття рішень [Електронний ресурс] – Режим доступу: http://uk.wikipedia.org/wiki/Дерево_прийняття_рішень
20. Загидиев А.М. Киберугрозі в банковской сфере // Научное сообщество студентов XXI столетия. Экономика: сб. ст. по мат. XXXI междунар. студ. науч.-практ. конф. № 4(31)

21. Інтелектуальний аналіз даних [Електронний ресурс]. – Режим доступу : <http://mirznanii.com/a/308854/ntelektualniy-analz-danikh>
22. Інтелектуальний аналіз даних [Електронний ресурс]. – Режим доступу : <https://univerfiles.com/1168907/Інтелектуальний-аналіз-даних/>
23. Киберугрозы и способы защиты финансовой безопасности 2017 [Електронний ресурс]. – Режим доступу : <https://ria-in.ru/it-industriya/kiberugrozy-i-sposoby-zashchity-finansovoj-bezopasnosti-2017>
24. Кібербезпека в банківській сфері [Електронний ресурс]. – Режим доступу : <https://icf.ua/blog/view/kiberbezopasnost-v-bankovskoy-sfere>
25. Кластерний аналіз [Електронний ресурс] – Режим доступу: http://uk.wikipedia.org/wiki/Кластерний_аналіз
26. Колодчак О.М. Інтелектуальний аналіз даних / О.М. Колодчак. – Національний університет «Львівська політехніка», 2013
27. Логістична регресія [Електронний ресурс] – Режим доступу: <http://www.statistica.ru/theory/logisticheskaya-regressiya/>
28. Марченко О.О., Россада Т.В. Актуальні проблеми Data Mining: Навчальний посібник для студентів факультету комп'ютерних наук та кібернетики. — Київ. — 2017. — 150 с.
29. Методи інтелектуального аналізу даних [Електронний ресурс]. – Режим доступу : <http://buklib.net/books/24506/>
30. Методы интеллектуального анализа данных [Електронний ресурс]. – Режим доступу : https://studme.org/1228112810027/ekonomika/metody_intellektualnogo_analiza_dannyh
31. Моделювання економіки : методичні рекомендації щодо виконання курсової роботи / Державний вищий навчальний заклад «Українська академія банківської справи Національного банку України» ; [уклад.: Л. П. Перхун, О. В. Кузьменко, С. М. Братушка та ін.]. – Суми : ДВНЗ «УАБС НБУ», 2014. – 34 с.

32. Нейронні мережі [Електронний ресурс]. – Режим доступу : <http://ukrbukva.net/page,2,92840-Neironnye-seti.html>
33. Нейронные сети [Електронний ресурс]. – Режим доступу : https://studopedia.su/14_173071_neyronnie-seti.html
34. Паклін Н. Логістична регресія та ROC-аналіз – математичний апарат [Електронний ресурс] / Н. Паклін – Режим доступу: <http://www.basegroup.ru/library/analysis/regression/logistic/>
35. Постанова НБУ Про затвердження Положення про організацію заходів із забезпечення інформаційної безпеки в банківській системі України [Електронний ресурс]. – Режим доступу : <https://bank.gov.ua/document/download?docId=56426049>
36. Проект Стратегії забезпечення кібернетичної безпеки України [Електронний ресурс]. – Режим доступу : http://www.niss.gov.ua/public/File/2013_nauk_an_rozrobku/kiberstrateg.pdf
37. Регрессионный анализ [Електронний ресурс]. – Режим доступу : http://57705.selcdn.com/MSU_2013/EM4.pdf
38. Теорія і практика розвитку наукових знань (частина II): матеріали II Міжнародної науково-практичної конференції м. Київ, 28-29 грудня 2017 року. – Київ.: МЦНД, 2017. – 56 с.
39. Чернышова Г.Ю. Интеллектуальный анализ данных: учебное пособие для студентов / Г.Ю.Чернышова. – Саратов : Саратовский государственный социально-экономический университет, 2012. – 92 с.
40. Штучна нейронна мережа [Електронний ресурс] – Режим доступу: http://uk.wikipedia.org/wiki/Штучна_нейронна_мережа

ДОДАТОК А

SUMMARY

Skovronska A. I. Intellectual analysis of cyber threats in banks. – Masters-level Qualification Thesis. Sumy State University, Sumy, 2018.

In this work are investigated the theoretical and methodological foundations of the intellectual analysis of cyber threats in banks. The analysis of existing approaches to detecting cyber threats in banks and methods of intellectual analysis is carried out. The main purpose of this study is to build a mathematical model for detecting cyber threats in banks and its practical implementation using intelligent analysis methods with the help of the package «SAS Enterprise Miner».

Keywords: intellectual analysis, cyberthreat, regression, decision tree, neural network.

АНОТАЦІЯ

Сковронська А. І. Інтелектуальний аналіз кіберзагроз в банках. – Кваліфікаційна магістерська робота. Сумський державний університет, Суми, 2018 р.

У роботі досліджено теоретико-методологічні основи інтелектуального аналізу кіберзагроз в банках. Проведений аналіз існуючих підходів до виявлення кіберзагроз у банках та методів інтелектуального аналізу. Основною метою цього дослідження є побудова математичної моделі для виявлення кіберзагроз у банках та її практична реалізація із використанням інтелектуального аналізу за допомогою пакету «SAS Enterprise Miner».

Ключові слова: інтелектуальний аналіз, кіберзагроза, регресія, дерево рішень, нейронна мережа.